



# Error estimates for the gradient discretisation of degenerate parabolic equation of porous medium type

Clément Cancès, Jerome Droniou, Cindy Guichard, Gianmarco Manzini,  
Manuela Bastidas Olivares, Iuliu Sorin Pop

## ► To cite this version:

Clément Cancès, Jerome Droniou, Cindy Guichard, Gianmarco Manzini, Manuela Bastidas Olivares, et al.. Error estimates for the gradient discretisation of degenerate parabolic equation of porous medium type. Daniele A. Di Pietro; Roland Masson; Luca Formaggia. Polyhedral methods in geosciences, Springer, In press, SEMA-SIMAI. hal-02540067

**HAL Id: hal-02540067**

**<https://hal.science/hal-02540067>**

Submitted on 10 Apr 2020

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# ERROR ESTIMATES FOR THE GRADIENT DISCRETISATION OF DEGENERATE PARABOLIC EQUATION OF POROUS MEDIUM TYPE

CLÉMENT CANCÈS, JÉRÔME DRONIOU, CINDY GUICHARD, GIANMARCO MANZINI,  
MANUELA BASTIDAS OLIVARES, AND IULIU SORIN POP

ABSTRACT. The gradient discretisation method (GDM) is a generic framework for the spatial discretisation of partial differential equations. The goal of this contribution is to establish an error estimate for a class of degenerate parabolic problems, obtained under very mild regularity assumptions on the exact solution. Our study covers well-known models like the porous medium equation and the fast diffusion equations, as well as the strongly degenerate Stefan problem. Several schemes are then compared in a last section devoted to numerical results.

## 1. INTRODUCTION

Degenerate parabolic equations appear as mathematical models for numerous real-life applications, like reactive solute transport in porous media, water infiltration in the vadose zone, geological CO<sub>2</sub> sequestration, oil recovery, biological systems, or phase transition problems. In the simplest form, one has

$$(1) \quad \begin{aligned} \partial_t \bar{u} - \Delta \zeta(\bar{u}) &= f && \text{in } (0, T) \times \Omega, \\ \zeta(\bar{u}) &= 0 && \text{on } (0, T) \times \partial\Omega, \\ \bar{u}(0, \cdot) &= u_{\text{ini}} && \text{on } \Omega. \end{aligned}$$

With  $L^\infty$  denoting the space of essentially bounded functions and  $\|\cdot\|_\infty$  the corresponding norm, throughout this chapter we assume the following.

(A0)  $T > 0$  and  $\Omega$  is a bounded connected domain of  $\mathbb{R}^d$  ( $d \in \mathbb{N}^*$ ) with Lipschitz continuous boundary  $\partial\Omega$ .

(A1)  $\zeta : \mathbb{R} \rightarrow \mathbb{R}$  is continuous, non-decreasing and satisfies  $\zeta(0) = 0$ .

(A2)  $u_{\text{ini}} \in L^\infty(\Omega)$  with  $M_0 := \|u_{\text{ini}}\|_\infty$ .

(A3)  $f \in L^\infty((0, T] \times \Omega)$  with  $M_f := \|f\|_\infty$ .

As follows from (A1),  $\zeta'$  may become zero, or unbounded for certain arguments  $\bar{u}$ . Consequently, the equation may degenerate from a parabolic equation into an elliptic or an ordinary one. The degeneracy regions are not known a-priori, but depend on the solution itself and may change in time.

One of the most representative example in this sense, the porous medium equation (PME), appeared in the last century as a mathematical model for the flow of an ideal gas in a porous medium. In this case one has

$$(2) \quad \zeta(\bar{u}) = |\bar{u}|^{m-1} \bar{u} \text{ for some } m > 1.$$

Compared to the heat equation, which is obtained for  $m = 1$  and in which the equation is linear and parabolic everywhere regardless of the data, if  $m > 1$  the nonlinear diffusive term vanishes if  $\bar{u} \leq 0$ , and the equation degenerates. In particular, this leads to the occurrence of free boundaries separating regions in  $\Omega$  where  $\bar{u} > 0$  from those where  $\bar{u} \leq 0$ . These free boundaries have an a-priori unknown location and move in time with a finite speed, which is the reason for calling such cases as "slow diffusion" ones.

Another remarkable example in the category of "slow diffusion" equations is the Stefan problem, which models phase transition problems like melting or solidification. In this case

$$(3) \quad \zeta(\bar{u}) = \begin{cases} \bar{u}, & \text{if } u < 0, \\ \max\{0, \bar{u} - 1\}, & \text{if } u \geq 1. \end{cases}$$

Though bounded,  $\zeta'$  is vanishing on the entire interval  $(0, 1)$ .

A different situation appears when  $\zeta$  is as in (2), but with  $m \in (0, 1)$ . In this case no free boundaries occur, but  $\zeta' \rightarrow \infty$  whenever  $\bar{u} \rightarrow 0$  so the diffusion coefficient becomes unbounded.

This equation is also known as generalized porous medium equation (GPME), and one speaks about a "fast diffusion". It can appear as a mathematical model for reactive transport in porous media, for equilibrium kinetics (see [6]).

The degeneracy has direct impact on the regularity of the solutions. Unlike the regular parabolic case, the solutions to degenerate parabolic problems have lower regularity, and the singularities are not smoothed out but may even develop in time. Such effects are particularly encountered at the free boundaries. The lack of regularity motivates the introduction of a notion of weak solution.

We use standard notations and function spaces in the functional analysis:  $L^2(\Omega)$ ,  $L^\infty(\Omega)$ ,  $H_0^1(\Omega)$ , or its dual  $H^{-1}(\Omega)$ . Whenever obvious, the domain  $\Omega$  is left out. With  $X$  being one of the spaces before,  $L^2(0, T; X)$  is the space of  $X$ -valued measurable functions that are square integrable in the sense of Bochner. We let  $(\cdot, \cdot)$  stand for the inner product on  $L^2(\Omega)$ , or the duality pairing between  $H_0^1(\Omega)$  and  $H^{-1}(\Omega)$ , and  $\|\cdot\|$  for the norm in  $L^2(\Omega)$ , or the straightforward extension to  $L^2(\Omega)^d$ , and  $\|\cdot\|_\infty$  is the  $L^\infty$  norm in  $\Omega$  or in  $(0, T] \times \Omega$ . We will often write  $u$  or  $u(t)$  instead of  $u(t, \mathbf{x})$  and use  $C$  to denote a generic positive constant independent of the discretisation parameters or the function itself.

We start by defining a weak solution for (1):

**Definition 1.1.** *A weak solution to (1) is a function  $\bar{u} \in H^1(0, T; H^{-1}(\Omega))$  s.t.  $\zeta(\bar{u}) \in L^2(0, T; H_0^1(\Omega))$ ,  $\bar{u}(0) = u_{\text{ini}}$  in  $H^{-1}(\Omega)$  and for a.e.  $t \in (0, T]$  and  $v \in H_0^1(\Omega)$  it holds*

$$(4) \quad (\partial_t \bar{u}(t), v) + (\nabla \zeta(\bar{u}(t)), \nabla v) = (f(t), v).$$

The existence and uniqueness of a weak solution to (1) is proved e.g. in [2] and [42] in the case where  $\zeta$  is increasing. If  $\zeta$  is merely nondecreasing, existence and uniqueness still hold, see e.g. [15], as well as [44]. As already suggested, the degenerate aspect of the problem makes the usual regularity theory for parabolic problems (see for instance [35]) fail. What is kept is mainly the following:

- *Maximum principle:* the solution  $\bar{u}$  belongs to  $L^\infty((0, T] \times \Omega)$ , with

$$(5) \quad \|\bar{u}\|_\infty \leq M_0 + TM_f.$$

- *Energy estimate:* Consider the primitive of  $\zeta$  defined by  $\Xi : \mathbb{R} \rightarrow \mathbb{R}$ ,  $\Xi(v) = \int_0^v \zeta(z) dz$ .  $\Xi$  is convex and positive and one has

$$(6) \quad \int_\Omega \Xi(\bar{u}(t)) + \frac{1}{2} \int_0^t \int_\Omega |\nabla \zeta(s)|^2 \leq \int_\Omega \Xi(u_{\text{ini}}) + \frac{1}{2} \|f\|_{L^2(0, T; H^{-1}(\Omega))}^2.$$

- *Continuity of  $\zeta(\bar{u})$ :* it is shown in [54] under quite general assumptions on  $\zeta$  (including cases where  $\zeta$  is constant on an interval) that  $\zeta(\bar{u})$  belongs to  $C((0, T] \times \Omega)$ . In the case where  $\zeta$  is increasing (thus invertible), one gets that  $\bar{u} \in C((0, T] \times \Omega)$  too. Because of the degeneracy of the problem, this estimate is not enough to initiate a bootstrap to recover the usual parabolic regularity theory.
- *Time continuity of  $\bar{u}$ :* even if  $\zeta$  is not invertible, one can still give a (weaker) sense to  $\bar{u}(t)$  as a function (and not only as a distribution in  $H^{-1}$  as suggested by Definition 1.1). Indeed,  $\bar{u} \in C([0, T]; L^p(\Omega))$  for all  $p \in [1, +\infty)$  thanks to [11].

Further regularity results in the PME case where  $\zeta(u) = |u|^{m-1}u$  (or more generally when  $\zeta$  is increasing) can be found in the monographs [49, 50] (see also [36] for the local Hölder continuity), while the Stefan problem is extensively discussed in [39].

The literature on the numerical approximation of degenerate parabolic equations is extremely rich. Often, the numerical scheme are combining a regularisation step, that is used to deal with the lack of regularity of the solution to degenerate problems. Whenever regularisation is involved, this is mostly obtained through a perturbation  $\zeta_\varepsilon$  of  $\zeta$ , of which derivative is bounded away from 0 from below and from infinity from above (see e.g. [41]). Alternatively, one can exploit the maximum principle and perturb the boundary and initial data in such a way that the solution stays away from values at which degeneracy is encountered.

Concerning various specific numerical schemes, we mention that often the time stepping is of first order. In particular Euler implicit of semi-implicit methods are popular, and this is due to the lack in regularity of the solution. For the spatial discretisation we mention the conformal finite element schemes analysed e.g. in [41] for the slow diffusion, or in [6] for the fast diffusion. The convergence of the mixed finite element discretisation is proved in [5, 52] for the slow diffusion case, and for

a range allowing for both kind of degeneracies in [48]. We also mention [53] for the analysis of a scheme combining mortars with mixed finite elements. These papers are proving the convergence of the scheme by obtaining a-priori error estimates rigorously. The convergence of finite volume schemes is proved in [3, 4, 29, 31] by means of compactness arguments, and in [30] for a finite volume phase-by-phase upstream weighting. Error estimates are obtained in [34] for a multipoint flux approximation scheme by using the equivalence with a mixed finite element scheme, and in [45] for the simplest two-point approximation in the slow diffusion case, but under minimal regularity assumptions. Discontinuous Galerkin schemes for porous media flow models leading to degenerate parabolic equations are analysed e.g. in [26, 27]. To conclude this paragraph, we mention that a posteriori error estimates for degenerate problems related to porous media flows are derived in [14], [51].

The goal of this chapter is to study in a general way a large class of numerical approximation of (1), entering the framework of the so-called *Gradient Discretisation Method (GDM)* [23]. This general framework is detailed in Section 3. The ideas used in the numerical analysis below apply to methods which are energetically stable (i.e., a discrete counterparts of (6) holds). Our approach does not require any monotonicity properties on the approximation, like the maximum principle. This choice is due to the fact that proving the maximum principle (5), as well as the (time-)continuity for  $\zeta(\bar{u})$  and  $\bar{u}$ , relies strongly on the monotonicity of equation (1), which also extends to the time-discrete case but not to the fully discrete provided by the GDM. We mention that the convergence of the GDM for a class of problems covering (1) is obtained in [28, 21] by means of compactness arguments. The aim here is to extend such results by providing a-priori error estimates

**Remark 1.2.** *The present results can be adapted without any particular further difficulty to the case of problems with anisotropy of the form*

$$\partial_t \bar{u} - \nabla \cdot (\Lambda \nabla \zeta(\bar{u})) = f \quad \text{in } (0, T) \times \Omega,$$

where  $\Lambda \in L^\infty(\Omega; \mathbb{R}^{d \times d})$  is a symmetric definite positive tensor field, i.e.,  $\Lambda(\mathbf{x}) = \Lambda(\mathbf{x})^T$  and there exists  $\lambda_m, \lambda_M > 0$  such that

$$\lambda_m |v|^2 \leq \Lambda(\mathbf{x}) v \cdot v \leq \lambda_M |v|^2, \quad \mathbf{x} \in \Omega, v \in \mathbb{R}^d.$$

The paper is organised as follows. Section 2 is introducing the sequence of semi-discrete in time problems. The time discretisation relies on the backward Euler scheme and is thus very standard. The a-priori error estimates for the time discretisation are deeply inspired from [38], and do not require any regularity assumption on the exact solution. Section 3 is devoted to the fully discrete setting. This encompasses the definition of the notions of *Gradient Discretisation* and *Gradient Scheme*, which were introduced in [25] and further developed in the monograph [23]. The main result is an error estimate for any scheme entering this general framework of the GDM. To this purpose, reasonable extra regularity slightly overpassing the aforementioned regularity results rigorously established in the literature will be assumed on the solution  $\bar{u}$ . Finally, several numerical schemes are compared in Section 4; these schemes consist of the Locally Enriched Non-Conforming Polytopal scheme, the Hybrid Mimetic Mixed method, two versions of the Vertex Approximate Gradient scheme, the mass-lumped  $\mathbb{P}^1$  Finite Element scheme, the Hybridizable Discontinuous Galerkin scheme, and the Conforming Virtual Element Method.

## 2. SEMI-DISCRETE IN TIME PROBLEM

Our purpose in this section is to show how to derive an error estimate using only minimal regularity assumptions for time-discrete approximations of (1). To this end, we first establish some a-priori estimates on the time-discrete solution. This section shall be seen as a first step towards the derivation of the fully discrete error estimate of Theorem 3.4.

**2.1. The time discretisation.** In view of the low regularity of the solution, we only consider first order time discretisation schemes. To this aim we consider a sequence of times  $0 = t^{(0)} < t^{(1)} < \dots < t^{(N)} = T$  ( $N \in \mathbb{N}^*$ ) and define the time steps  $\delta t^{(n+\frac{1}{2})} = t^{(n+1)} - t^{(n)}$  ( $n \in \{0, \dots, N-1\}$ ). We let  $\bar{u}^{(n)}$  be a time discrete approximation of  $\bar{u}(t^{(n)})$ . To define a weak solution to the time-discrete problems we use the set  $X_0 := \{\Xi(u) \in L^1(\Omega) : \zeta(u) \in H_0^1(\Omega)\}$ . The Euler implicit discretisation of (4) consists in finding a sequence of solutions to the time discrete problems, as defined in

**Definition 2.1** (Time semi-discrete problem). *Set  $\bar{u}^{(0)} = u_{\text{ini}}$ . With  $n \in \{0, \dots, N-1\}$ , given  $\bar{u}^{(n)} \in X_0$ , a weak solution  $\bar{u}^{(n+1)} \in X_0$  to the time discrete problem at time step  $t^{(n+1)}$  satisfies, for all  $v \in H_0^1(\Omega)$ ,*

$$(7) \quad (\bar{u}^{(n+1)}, v) + \delta t^{(n+\frac{1}{2})} (\nabla \zeta(\bar{u}^{(n+1)}), \nabla v) = (\bar{u}^{(n)}, v) + \delta t^{(n+\frac{1}{2})} (f^{(n+1)}, v),$$

where  $f^{(n+1)}(\mathbf{x}) = \frac{1}{\delta t^{(n+\frac{1}{2})}} \int_{t^{(n)}}^{t^{(n+1)}} f(s, \mathbf{x}) ds$ .

**Theorem 2.2** (Existence and uniqueness of a solution to the semi-discrete problem). *There exists a unique family  $(\bar{u}^{(n)})_{n=0, \dots, N}$  solution to the semi-discrete problem in the sense of Definition 2.1.*

*Proof.* Follows by applying [22, Theorem A.1] to solve, at each step, the non-linear elliptic problem  $w - \delta t^{(n+\frac{1}{2})} \Delta \zeta(w) = \bar{u}^{(n)} + \delta t^{(n+\frac{1}{2})} f^{(n+1)}$ .  $\square$

**2.2. A-priori estimates.** Our goal is to provide a fully discrete error analysis for numerical schemes for (1). For ease of legibility, we start by discussing some properties of the semi-discrete, Euler implicit discretisation in (7). In doing so, we follow the ideas in [38], where the convergence of a linear, time discrete scheme is proved for a class of problems that includes (1).

We start with a remark on the essential boundedness of a solution. This property is physically justified for many of the applications that can be modelled mathematically in the form of (1) (e.g. the gas flow in porous media flows, or the reactive transport). In this context, the essential boundedness is inherited by the time discrete solutions, which satisfy a maximum principle. However, since this property does not extend to the fully discrete cases excepting some particular finite element or finite volume discretisation, we will avoid using it below.

Assuming that the initial data and the source term are both essentially bounded, as stated in Assumptions (A2) and (A3), one has

**Lemma 2.3.** *Assume  $\bar{u}^{(n)} \in X_0$  is such that  $\|\bar{u}^{(n)}\|_\infty \leq M_0 + M_f t^{(n)}$ . Then the solution  $\bar{u}^{(n+1)}$  of (7) satisfies  $\|\bar{u}^{(n+1)}\|_\infty \leq M_0 + M_f t^{(n+1)}$ .*

These estimates are obtained straightforwardly by testing in (7) with  $v = [\zeta(\bar{u}^{(n+1)}) - \zeta(M_0 + M_f t^{(n+1)})]_+$ , and with  $v = [\zeta(\bar{u}^{(n+1)}) + \zeta(M_0 + M_f t^{(n+1)})]_-$  (with  $[s]_+ = \max(s, 0)$  and  $[s]_- = \min(s, 0)$ ). We omit the details.

We state some elementary results that are used below, and which are valid for all set of vectors  $\mathbf{a}_n, \mathbf{b}_n \in \mathbb{R}^d$  ( $d \geq 1$ ),  $n \in \{0, \dots, m\}$ .

$$(8) \quad 2 \sum_{n=1}^m \mathbf{a}_n \cdot (\mathbf{a}_n - \mathbf{a}_{n-1}) = |\mathbf{a}_m|^2 - |\mathbf{a}_0|^2 + \sum_{n=1}^m |\mathbf{a}_n - \mathbf{a}_{n-1}|^2,$$

$$(9) \quad 2 \sum_{n=0}^m \sum_{j=0}^n \mathbf{a}_n \cdot \mathbf{a}_j = \left| \sum_{n=0}^m \mathbf{a}_n \right|^2 + \sum_{n=0}^m |\mathbf{a}_n|^2,$$

$$(10) \quad \sum_{n=1}^m \mathbf{a}_n \cdot (\mathbf{b}_n - \mathbf{b}_{n-1}) = \mathbf{a}_m \cdot \mathbf{b}_m - \mathbf{a}_0 \cdot \mathbf{b}_0 - \sum_{n=1}^m (\mathbf{a}_n - \mathbf{a}_{n-1}) \cdot \mathbf{b}_{n-1}.$$

Further, with the convex, positive primitive  $\Xi$  of  $\zeta$  appearing in (6), a classical convexity relation yields

$$(11) \quad (b - a)\zeta(b) \geq \Xi(b) - \Xi(a), \quad \forall a, b \in \mathbb{R}.$$

Finally, we state for completeness the Young inequality, valid for any  $a, b \in \mathbb{R}$  and  $\varepsilon > 0$ ,

$$(12) \quad ab \leq \frac{1}{2\varepsilon} a^2 + \frac{\varepsilon}{2} b^2.$$

The stability of the time discrete scheme is stated in

**Lemma 2.4.** *Let  $(\bar{u}^{(n+1)})_{n=0, \dots, N-1}$  be the sequence of time discrete solutions introduced in Definition 2.1. Then,*

$$(13) \quad \sum_{n=0}^{N-1} \delta t^{(n+\frac{1}{2})} \|\nabla \zeta(\bar{u}^{(n+1)})\|^2 \leq C.$$

*Proof.* Taking in (7)  $v = \zeta(\bar{u}^{(n+1)})$  gives

$$(14) \quad (\bar{u}^{(n+1)} - \bar{u}^{(n)}, \zeta(\bar{u}^{(n+1)})) + \delta t^{(n+\frac{1}{2})} \|\nabla \zeta(\bar{u}^{(n+1)})\|^2 = \delta t^{(n+\frac{1}{2})} (f^{(n+1)}, \zeta(\bar{u}^{(n+1)})).$$

For the first term one uses (11) to obtain

$$(15) \quad (\bar{u}^{(n+1)} - \bar{u}^{(n)}, \zeta(\bar{u}^{(n+1)})) \geq \int_{\Omega} \Xi(\bar{u}^{(n+1)}) - \Xi(\bar{u}^{(n)}) \, dx.$$

The second term needs no further discussion, whereas for the term on the right one obtains

$$(16) \quad \delta t^{(n+\frac{1}{2})} |(f^{(n+1)}, \zeta(\bar{u}^{(n+1)}))| \leq \frac{\delta t^{(n+\frac{1}{2})}}{2} \|f^{(n+1)}\|_{H^{-1}(\Omega)}^2 + \frac{\delta t^{(n+\frac{1}{2})}}{2} \|\nabla \zeta(\bar{u}^{(n+1)})\|^2.$$

Using (15) and (16) into (14) and summing the resulting relation over  $n \in \{1, \dots, N-1\}$  yields

$$(17) \quad \int_{\Omega} \Xi(\bar{u}^{(n+1)}) + \frac{1}{2} \sum_{n=0}^{N-1} \delta t^{(n+\frac{1}{2})} \|\nabla \zeta(\bar{u}^{(n+1)})\|^2 \leq \int_{\Omega} \Xi(\bar{u}^{(0)}) + \sum_{n=0}^{N-1} \frac{\delta t^{(n+\frac{1}{2})}}{2} \|f^{(n+1)}\|_{H^{-1}(\Omega)}^2.$$

The first term on the right is bounded due to Assumption (A2). For the second term on the right one uses Assumption (A3) and the fact that  $L^\infty((0, T) \times \Omega)$  is continuously embedded into  $L^2(0, T; H^{-1}(\Omega))$ , to obtain a uniform bound w.r.t.  $N$ . Since  $\Xi$  is a positive function, the first term on the left is positive, which provides (13).  $\square$

**Remark 2.5** (Other estimates). *Other a-priori estimates can be obtained if further assumptions are made on  $\zeta$  and on the initial data. For example, if  $\zeta$  is Lipschitz and  $\zeta(u_{\text{ini}}) \in H_0^1(\Omega)$ , then one can prove that  $\|\nabla \zeta(\bar{u})\|$  and  $\|\nabla \zeta(\bar{u}^{(n)})\|$  are uniformly bounded (w.r.t.  $t$ , respectively  $n$ ) and obtain  $L^2$  estimates for  $\partial_t \zeta(\bar{u})$  or its time discrete counterpart.*

**2.3. Error estimates, time discrete case.** We now establish error estimates for the time discrete scheme. The proof follows the lines of [38, Section 4]. We use the following notations for the errors

$$\begin{aligned} e_u(t) &:= \bar{u}(t) - \bar{u}^{(n+1)}, \\ e_\zeta(t) &:= \zeta(\bar{u}(t)) - \zeta(\bar{u}^{(n+1)}), \\ e_\gamma(t) &:= \gamma(\bar{u}(t)) - \gamma(\bar{u}^{(n+1)}), \end{aligned}$$

for  $t \in (t^{(n)}, t^{(n+1)}]$  and  $n \in \{0, \dots, N-1\}$ . With this, one has

**Theorem 2.6.** *Let  $\bar{u}$  be the solution in Definition 1.1, and  $(\bar{u}^{(n+1)})_{n=0, \dots, N-1}$  be the sequence of solutions to the time discrete problems in Definition 2.1. Setting  $\delta t = \max_{n \in \{0, \dots, N-1\}} \delta t^{(n+\frac{1}{2})}$ , one has*

$$(18) \quad \max_{n \in \{0, \dots, N-1\}} \|e_u(t^{(n+1)})\|_{H^{-1}(\Omega)}^2 + \int_0^T (e_u(t), e_\zeta(t)) \, dt \leq C \delta t.$$

**Remark 2.7.** *Although the second term in (18) is not a proper norm, it can generate an error estimate in a certain norm whenever the particular form of  $\zeta$  is taken into account. This idea is exploited, in the fully discrete setting, in Corollary 3.7 below. Moreover, as for the a-priori estimates, under additional assumptions on  $\zeta$  one can get  $L^2$  error estimates for either  $\zeta(\bar{u})$  (if  $\zeta$  is Lipschitz continuous, as appearing in the slow diffusion case) or  $\bar{u}$  (if  $\zeta$  is bijective and its inverse Lipschitz, as appearing in the fast diffusion case).*

*Proof.* Before giving the proof we observe that, due to the monotonicity of  $\zeta$ , the second term in (18) is positive.

With  $j \in \{0, \dots, N-1\}$  we integrate (4) in time for  $t \in (t^{(j)}, t^{(j+1)}]$  and subtract (7) for  $n = j$  to obtain

$$(e_u(t^{(j+1)}) - e_u(t^{(j)}), v) + \left( \nabla \int_{t^{(j)}}^{t^{(j+1)}} e_\zeta(t) \, dt, \nabla v \right) = 0,$$

for all  $v \in H_0^1(\Omega)$ . After summation over  $j \in \{0, \dots, p\}$  for some  $p \in \{0, \dots, N-1\}$ , this yields

$$(19) \quad (e_u(t^{(p+1)}), v) + \left( \nabla \int_0^{t^{(p+1)}} e_\zeta(t) \, dt, \nabla v \right) = 0,$$

for all  $v \in H_0^1(\Omega)$ . Taking  $v = \int_{t^{(p)}}^{t^{(p+1)}} e_\zeta(t) dt$  in the above equation provides

$$(20) \quad \left( e_u(t^{(p+1)}), \int_{t^{(p)}}^{t^{(p+1)}} e_\zeta(t) dt \right) + \sum_{j=0}^p \left( \nabla \int_{t^{(j)}}^{t^{(j+1)}} e_\zeta(t) dt, \nabla \int_{t^{(p)}}^{t^{(p+1)}} e_\zeta(t) dt \right) = 0.$$

Fixing  $n \in \{0, \dots, N-1\}$ , we sum (20) over  $p \in \{0, \dots, n\}$  and obtain

$$(21) \quad \underbrace{\sum_{p=0}^n \left( e_u(t^{(p+1)}), \int_{t^{(p)}}^{t^{(p+1)}} e_\zeta(t) dt \right)}_{=: I_1} + \underbrace{\sum_{p=0}^n \sum_{j=0}^p \left( \nabla \int_{t^{(j)}}^{t^{(j+1)}} e_\zeta(t) dt, \nabla \int_{t^{(p)}}^{t^{(p+1)}} e_\zeta(t) dt \right)}_{=: I_2} = 0.$$

The first term can be rewritten as

$$(22) \quad I_1 = \underbrace{\int_0^{t^{(n+1)}} (e_u(t), e_\zeta(t)) dt}_{=: I_{11}} + \underbrace{\sum_{p=0}^n \int_{t^{(p)}}^{t^{(p+1)}} (\bar{u}(t^{(p+1)}) - \bar{u}(t), e_\zeta(t)) dt}_{=: I_{12}}.$$

Being positive,  $I_{11}$  needs no further handling. For  $I_{12}$  we write

$$\bar{u}(t^{(p+1)}) - \bar{u}(t) = \int_t^{t^{(p+1)}} \partial_s \bar{u} ds$$

to obtain

$$\begin{aligned} |I_{12}| &= \left| \sum_{p=0}^n \int_{t^{(p)}}^{t^{(p+1)}} \left( \int_t^{t^{(p+1)}} \partial_s \bar{u} ds, e_\zeta(t) \right) dt \right| \\ &\leq \left( \max_{n \in \{0, \dots, N-1\}} \delta t^{(n+\frac{1}{2})} \right) \|\partial_t \bar{u}\|_{L^2(0,T;H^{-1}(\Omega))} \|\nabla e_\zeta\|_{L^2(0,T;L^2(\Omega))}. \end{aligned}$$

The regularity of the weak solution prescribed in Definition 1.1 and the a-priori estimate (13) ensure the existence of a  $C > 0$  not depending on the time discretisation so that

$$\|\partial_t \bar{u}\|_{L^2(0,T;H^{-1}(\Omega))} \leq C, \quad \|\nabla e_\zeta\|_{L^2(0,T;L^2(\Omega))} \leq C.$$

As a consequence, we obtain that

$$(23) \quad |I_{12}| \leq C \delta t.$$

Finally, using (9),  $I_2$  is nonnegative and can be underestimated by

$$(24) \quad I_2 \geq \frac{1}{2} \left\| \nabla \int_0^{t^{(n+1)}} e_\zeta(t) dt \right\|^2.$$

Since  $n$  was chosen arbitrarily, using (22)–(24) in (21) yields

$$(25) \quad \int_0^T (e_u(t), e_\zeta(t)) dt + \max_{n \in \{0, \dots, N-1\}} \left\| \nabla \int_0^{t^{(n+1)}} e_\zeta(t) dt \right\|^2 \leq C \delta t.$$

To complete the proof of Theorem 2.6 one needs to estimate  $\|e_u\|_{H^{-1}(\Omega)}$ . This follows straightforwardly from (19). For all  $v \in H_0^1(\Omega)$  such that  $\|v\|_{H_0^1(\Omega)} \leq 1$ , the Cauchy–Schwarz inequality yields

$$(e_u(t^{(p+1)}), v) = - \left( \nabla \int_0^{t^{(p+1)}} e_\zeta(t) dt, \nabla v \right) \leq \left\| \nabla \int_0^{t^{(p+1)}} e_\zeta(t) dt \right\|.$$

Taking the supremum over such  $v$ , squaring and using (25) we infer  $\|e_u(t^{(p+1)})\|_{H^{-1}(\Omega)}^2 \leq C \delta t$  and the proof is complete.  $\square$



## 3. GRADIENT DISCRETISATION METHOD AND GENERIC ERROR ESTIMATE

**3.1. Definition of the gradient scheme.** The principle of the GDM is to replace, in the weak formulation of the problem, the continuous space and differential operators by discrete ones. To this aim a discrete space and function/gradient reconstructions on this space are used. Altogether these form a gradient discretisation (GD)  $\mathcal{D}$ . In general, very few assumptions are made on the GD [23]. However, to deal with the non-linearity we will need, in a similar way as in [22], to consider *nodal* gradient discretisations with *piecewise constant* reconstructions, which also contain the definition of an interpolator. Therefore, we take  $\mathcal{D} = (X_{\mathcal{D},0}, \Pi_{\mathcal{D}}, \nabla_{\mathcal{D}}, I_{\mathcal{D}})$  with

- (*Space*)  $X_{\mathcal{D},0} = \{v = (v_i)_{i \in I} : v_i = 0 \quad \forall i \in I_{\partial}\}$ , where  $I$  is a finite set and  $I_{\partial} \subset I$  identifies the boundary degrees of freedom.
- (*Function reconstruction*) Given a partition  $(U_i)_{i \in I}$  of  $\Omega$ , for all  $v = (v_i)_{i \in I} \in X_{\mathcal{D},0}$ , the reconstructed function  $\Pi_{\mathcal{D}}v \in L^\infty(\Omega)$  is defined by

$$(26) \quad \Pi_{\mathcal{D}}v = \sum_{i \in I} v_i \mathbf{1}_{U_i},$$

where  $\mathbf{1}_{U_i}$  is the characteristic function of  $U_i$ .

- (*Gradient reconstruction*)  $\nabla_{\mathcal{D}} : X_{\mathcal{D},0} \rightarrow L^2(\Omega)^d$  is a linear operator such that  $\|\nabla_{\mathcal{D}} \cdot\|$  is a norm on  $X_{\mathcal{D},0}$ .
- (*Interpolator*) The unknowns represent values at points  $(\mathbf{x}_i)_{i \in \bar{\Omega}}$ , with  $\mathbf{x}_i \in \bar{U}_i$  for all  $i \in I$  and  $\mathbf{x}_i \in \partial\Omega$  whenever  $i \in I_{\partial}$ . With  $C_{\text{pw},0}(\Omega)$  the set of piecewise continuous functions on  $\bar{\Omega}$  that have a zero limit on  $\partial\Omega$ , we define the interpolator  $I_{\mathcal{D}} : C_{\text{pw},0}(\Omega) \rightarrow X_{\mathcal{D},0}$  such that

$$(27) \quad (I_{\mathcal{D}}\phi)_i = \text{ess-limsup}_{\mathbf{x} \rightarrow \mathbf{x}_i, \mathbf{x} \in \bar{U}_i} \phi(\mathbf{x}), \quad \forall i \in I,$$

where

$$\text{ess-limsup}_{\mathbf{x} \rightarrow \mathbf{x}_i, \mathbf{x} \in \bar{U}_i} \phi(\mathbf{x}) = \lim_{\epsilon \rightarrow 0} \text{ess-sup}_{B(\mathbf{x}_i, \epsilon) \cap \bar{U}_i} \phi.$$

The fact that the function reconstruction is piecewise constant enables us to define, for  $g : \mathbb{R} \rightarrow \mathbb{R}$  with  $g(0) = 0$  and  $v \in X_{\mathcal{D},0}$ , the element  $g(v) \in X_{\mathcal{D},0}$  by applying  $g$  to each nodal value: if  $v = (v_i)_{i \in I}$ , we set  $g(v) = (g(v_i))_{i \in I}$ . It then holds

$$(28) \quad \Pi_{\mathcal{D}}g(v) = g(\Pi_{\mathcal{D}}v), \quad \forall v \in X_{\mathcal{D},0}.$$

The subtle choice for the definition of the interpolator is motivated by the following points. Since our study covers the Stefan problem, whose solution might be discontinuous, there is a real need to define an interpolator that allows for merely piecewise continuous functions. If  $\phi$  is continuous, say  $\phi \in C_0(\Omega)$ , then

$$\text{ess-limsup}_{\mathbf{x} \rightarrow \mathbf{x}_i, \mathbf{x} \in \bar{U}_i} \phi(\mathbf{x}) = \phi(\mathbf{x}_i), \quad \forall i \in I.$$

Therefore, for any continuous function  $g : \mathbb{R} \rightarrow \mathbb{R}$  and  $\phi \in C_{\text{pw},0}(\Omega)$ ,

$$(29) \quad I_{\mathcal{D}}g(\phi) = g(I_{\mathcal{D}}\phi).$$

When  $\phi$  is only piecewise continuous, then (29) still holds as soon as  $g$  is continuous and nondecreasing.

We can now define the gradient scheme for (1) with implicit time stepping. It is obtained from (7) using the discrete space for trial and test function, and replacing the functions and gradients by the corresponding reconstructions. This gives a sequence of fully discrete, nonlinear algebraic problems, obtained for  $n \in \{0, \dots, N-1\}$  and starting with  $u^{(0)} = I_{\mathcal{D}}u_{\text{ini}}$ .

**Problem  $\mathcal{P}_{\mathcal{D}}^{(n+1)}$ :** Given  $u^{(n)} \in X_{\mathcal{D},0}$ , find  $u^{(n+1)} \in X_{\mathcal{D},0}$  such that

$$(30) \quad \int_{\Omega} \Pi_{\mathcal{D}}(u^{(n+1)} - u^{(n)}) \Pi_{\mathcal{D}}v + \delta t^{(n+\frac{1}{2})} \int_{\Omega} \nabla_{\mathcal{D}} \zeta(u^{(n+1)}) \cdot \nabla_{\mathcal{D}}v = \delta t^{(n+\frac{1}{2})} \int_{\Omega} f^{(n+1)} \Pi_{\mathcal{D}}v$$

for all  $v \in X_{\mathcal{D},0}$ , where  $f^{(n+1)}$  is introduced in Definition 2.1.

**Proposition 3.1** (Existence and uniqueness of a solution to the gradient scheme [22, Lemma 2.7]). *There exists a solution  $u = (u^{(n)})_{n=0, \dots, N}$  to the gradient scheme and, if  $u_1, u_2 \in X_{\mathcal{D},0}^{N+1}$  are two solutions to this scheme, then  $\zeta(u_1) = \zeta(u_2)$  and  $\Pi_{\mathcal{D}}u_1 = \Pi_{\mathcal{D}}u_2$ .*



**Remark 3.2** (Limit to the uniqueness). *If  $\zeta$  is strictly increasing, then we have complete uniqueness of the solution:  $u_1 = u_2$ . However, when  $\zeta$  has plateau this uniqueness may fail [22, Remark 2.8].*

The accuracy of a GD is measured through three quantities: a discrete Poincaré constant  $C_{\mathcal{D}}$  (yielding the coercivity of the method), a measure of the defect of the discrete Stokes formula  $W_{\mathcal{D}}$  (associated with the limit-conformity of the method), and a measure of the interpolation error  $S_{\mathcal{D}}$  (which, when it tends to zero, yields the consistency of the method). The discrete Poincaré constant is

$$(31) \quad C_{\mathcal{D}} = \max_{v \in X_{\mathcal{D},0} \setminus \{0\}} \frac{\|\Pi_{\mathcal{D}} v\|}{\|\nabla_{\mathcal{D}} v\|}.$$

The measure of the defect of the discrete Stokes formula is  $W_{\mathcal{D}} : H_{\text{div}}(\Omega) \rightarrow [0, \infty)$  where, for all  $\psi \in H_{\text{div}}(\Omega)$  (that is,  $\psi \in L^2(\Omega)^d$  and  $\text{div} \psi \in L^2(\Omega)$ ),

$$(32) \quad W_{\mathcal{D}}(\psi) := \max_{v \in X_{\mathcal{D},0} \setminus \{0\}} \frac{1}{\|\nabla_{\mathcal{D}} v\|} \left| \int_{\Omega} \Pi_{\mathcal{D}} v \text{div} \psi + \nabla_{\mathcal{D}} v \cdot \psi \right|.$$

In the GDM, the interpolation error usually involves  $L^2$ -error in both function and gradient approximation. However, for time-dependent problems such as (1), it will be more efficient to use a weaker norm for the function approximation. We define the discrete  $H^{-1}$ -seminorm by: for  $\phi \in L^2(\Omega)$ ,

$$(33) \quad |\phi|_{\mathcal{D},*} := \max \left\{ \int_{\Omega} \phi \Pi_{\mathcal{D}} v : v \in X_{\mathcal{D},0}, \|\nabla_{\mathcal{D}} v\| \leq 1 \right\}.$$

We then set, for  $\phi \in C_{\text{pw},0}(\Omega)$  and  $\psi \in C_{\text{pw},0}(\Omega) \cap H_0^1(\Omega)$ ,

$$(34) \quad S_{\mathcal{D}}^{\Pi,*}(\phi) = |\Pi_{\mathcal{D}} I_{\mathcal{D}} \phi - \phi|_{\mathcal{D},*} \quad \text{and} \quad S_{\mathcal{D}}^{\nabla}(\psi) = \|\nabla_{\mathcal{D}} I_{\mathcal{D}} \psi - \nabla \psi\|.$$

**3.2. A-priori estimates.** We start with the observation that some properties of the solution to the original problem (1), or its time discrete counterpart, are not preserved by the gradient scheme (30). In particular we refer to the maximum principle (see Lemma 2.3), which in the spatially-continuous case is obtained by testing with a cut-off function. However, in the spatially discrete case the cut-off of an element in the finite dimensional space may not belong to that space any more, since the transition from negative to positive values does not necessarily happen at edges or nodes. For obtaining a priori estimates we therefore restrict to using in (30) test functions that are affine functions of  $\zeta(u^{(n+1)})$ , but mention that schemes allowing to take nonlinear functions are designed and analysed in [10, 12, 13].

Specifically, we extend in this section the estimates in Lemma 2.4 to the fully discrete case.

**Lemma 3.3.** *For the sequence of fully discrete solutions of (30) it holds*

$$(35) \quad \sum_{n=0}^{N-1} \delta t^{(n+\frac{1}{2})} \|\nabla_{\mathcal{D}} \zeta(u^{(n+1)})\|^2 \leq C_{\mathcal{D}}^2 \|f\|_{L^2(0,T;L^2(\Omega))}^2 + 2 \|\Xi(\Pi_{\mathcal{D}} u^{(0)})\|_{L^1(\Omega)}.$$

*Proof.* Choosing  $v = \zeta(u^{(n+1)})$  in (30) leads to

$$(36) \quad \int_{\Omega} \Pi_{\mathcal{D}}(u^{(n+1)} - u^{(n)}) \Pi_{\mathcal{D}} \zeta(u^{(n+1)}) + \delta t^{(n+\frac{1}{2})} \int_{\Omega} \left| \nabla_{\mathcal{D}} \zeta(u^{(n+1)}) \right|^2 = \delta t^{(n+\frac{1}{2})} \int_{\Omega} f^{(n+1)} \Pi_{\mathcal{D}} \zeta(u^{(n+1)})$$

for all  $n \in \{0, \dots, N-1\}$ . Using (28) and the convexity inequality (11), we have

$$(37) \quad \begin{aligned} \int_{\Omega} \Pi_{\mathcal{D}}(u^{(n+1)} - u^{(n)}) \Pi_{\mathcal{D}} \zeta(u^{(n+1)}) &= \int_{\Omega} (\Pi_{\mathcal{D}} u^{(n+1)} - \Pi_{\mathcal{D}} u^{(n)}) \zeta(\Pi_{\mathcal{D}} u^{(n+1)}) \\ &\geq \int_{\Omega} (\Xi(\Pi_{\mathcal{D}} u^{(n+1)}) - \Xi(\Pi_{\mathcal{D}} u^{(n)})). \end{aligned}$$

The right-hand side of (36) can be estimated thanks to Young and discrete Poincaré inequalities as follows:

$$(38) \quad \begin{aligned} \int_{\Omega} f^{(n+1)} \Pi_{\mathcal{D}} \zeta(u^{(n+1)}) &\leq \frac{C_{\mathcal{D}}^2}{2} \|f^{(n+1)}\|^2 + \frac{1}{2C_{\mathcal{D}}^2} \|\Pi_{\mathcal{D}} \zeta(u^{(n+1)})\|^2 \\ &\leq \frac{C_{\mathcal{D}}^2}{2} \|f^{(n+1)}\|^2 + \frac{1}{2} \|\nabla_{\mathcal{D}} \zeta(u^{(n+1)})\|^2. \end{aligned}$$

Combining (37) and (38) in (36), summing over  $n \in \{0, \dots, N-1\}$ , and using  $\Xi \geq 0$ , the proof of (35) is complete.  $\square$

**3.3. Error estimate.** With  $I_{\mathcal{D}}^n \bar{u} = I_{\mathcal{D}} \bar{u}(t^{(n)})$  for  $n \in \{0, \dots, N\}$ , we define the errors in  $X_{\mathcal{D},0}$ :

$$\begin{aligned} e_{\mathcal{D},u}^{(n)} &:= u_{\mathcal{D}}^{(n)} - I_{\mathcal{D}}^n \bar{u}, \\ e_{\mathcal{D},\zeta}^{(n)} &:= \zeta(u_{\mathcal{D}}^{(n)}) - I_{\mathcal{D}}^n \zeta(\bar{u}), \end{aligned}$$

as well as, for  $n \in \{1, \dots, N\}$ ,

$$\varepsilon_{\mathcal{D},\zeta}^{(n)} := \sum_{p=1}^n \delta t^{(p-\frac{1}{2})} e_{\mathcal{D},\zeta}^{(p)}.$$

The error estimates for the fully discrete approximation is stated in the following theorem, whose proof is carried out in Section 3.4.

**Theorem 3.4** (GDM error estimate for degenerate parabolic problem). *Assume that the solution  $\bar{u}$  of (1) satisfies  $\bar{u}(t, \cdot) \in C_{\text{pw},0}(\Omega)$  for all  $t \in [0, T]$ ,  $\zeta(\bar{u}) \in C([0, T]; C_0(\Omega) \cap H_0^1(\Omega))$ , and  $\nabla \zeta(\bar{u}) \in C([0, T]; H_{\text{div}}(\Omega))$ . Then, there exists a universal constant  $K$  depending neither on the data of the continuous problem nor on the discretisation parameters such that*

$$\begin{aligned} (39) \quad \max_{1 \leq n \leq N} \left| \Pi_{\mathcal{D}} e_{\mathcal{D},u}^{(n)} \right|_{\mathcal{D},*}^2 + \max_{1 \leq n \leq N} \left\| \nabla_{\mathcal{D}} \varepsilon_{\mathcal{D},\zeta}^{(n)} \right\|^2 \\ + \sum_{n=0}^{N-1} \delta t^{(n+\frac{1}{2})} (\Pi_{\mathcal{D}} e_{\mathcal{D},u}^{(n+1)}, \Pi_{\mathcal{D}} e_{\mathcal{D},\zeta}^{(n+1)}) \leq K(1+T) E_{\mathcal{D}}(\bar{u})^2, \end{aligned}$$

where

$$(40) \quad E_{\mathcal{D}}(\bar{u})^2 = \sum_{n=0}^{N-1} \delta t^{(n+\frac{1}{2})} (E_{\mathcal{D}}^n(\bar{u}))^2$$

with, for  $n \in \{0, \dots, N-1\}$ ,

$$\begin{aligned} (41) \quad E_{\mathcal{D}}^n(\bar{u}) &:= \left| \frac{1}{\delta t^{(n+\frac{1}{2})}} \int_{t^{(n)}}^{t^{(n+1)}} \Delta \zeta(\bar{u}(s)) ds - \Delta \zeta(\bar{u}(t^{(n+1)})) \right|_{\mathcal{D},*} \\ &+ S_{\mathcal{D}}^{\Pi,*} \left( \frac{\bar{u}(t^{(n+1)}) - \bar{u}(t^{(n)})}{\delta t^{(n+\frac{1}{2})}} \right) + S_{\mathcal{D}}^{\nabla}(\zeta(\bar{u}(t^{(n+1)}))) \\ &+ W_{\mathcal{D}}(\nabla \zeta(\bar{u}(t^{(n+1)}))). \end{aligned}$$

Due to the non-decreasing property of  $\zeta$ , each term in the sum in the left-hand side of (39) is non-negative. For the nonlinearities appearing in the case of the Stefan or the porous medium equations, (39) leads to the following error estimates on more natural quantities.

**Remark 3.5** (Expected rates of convergence). *Under regularity assumptions on  $\bar{u}$ , a rate of convergence in terms of time and mesh sizes can be obtained on  $E_{\mathcal{D}}(\bar{u})$ . Specifically, if  $\zeta(\bar{u}) \in C([0, T]; H^2(\Omega))$ ,  $\partial_t \bar{u} \in L^\infty(0, T; H_0^1(\Omega))$  and  $\Delta \zeta(\bar{u}) \in W^{1,\infty}(0, T; L^2(\Omega))$ , following the techniques in [23, Section 7.4] and using  $|\phi|_{\mathcal{D},*} \leq C_{\mathcal{D}} \|\phi\|$  it can be proved, for all usual low-order gradient discretisations based on meshes of maximum size  $h$  (which include all schemes used in Section 4 except VAG-b), that  $E_{\mathcal{D}}^n(\bar{u}) \leq C_{\bar{u}}((1 + C_{\mathcal{D}})\delta t^{(n+\frac{1}{2})} + h)$ , where  $C_{\bar{u}}$  only depends on  $\bar{u}$ . Hence, in this situation and setting  $\delta t = \max_{n \in \{0, \dots, N-1\}} \delta t^{(n+\frac{1}{2})}$ , we have*

$$E_{\mathcal{D}}(\bar{u}) \leq T^{\frac{1}{2}} C_{\bar{u}}((1 + C_{\mathcal{D}})\delta t + h).$$

**Remark 3.6.** *For the slow diffusion case and under the regularity stated in Definition 1.1, error estimates for a simple, two-point flux approximation scheme (which fits in the GDM framework) are obtained in [45]. The approach there consists in using a discrete Green function to estimate the error for the fully discrete approximation of the sequence of time discrete approximations (Definition 2.1). This approach involves a regularisation step, which we avoid here by using a different strategy.*

**Corollary 3.7** (Estimate for the Stefan equation and the PME). *Under the assumptions and notations in Theorem 3.4, the following holds.*

- (Stefan equation) Assume that  $\zeta$  is Lipschitz-continuous with Lipschitz constant  $L_\zeta$ . Then,

$$(42) \quad \left( \sum_{n=0}^{N-1} \delta t^{(n+\frac{1}{2})} \|\zeta(\Pi_{\mathcal{D}} u^{(n+1)}) - \zeta(\Pi_{\mathcal{D}} I_{\mathcal{D}}^{n+1} \bar{u})\|^2 \right)^{\frac{1}{2}} \leq (K(1+T)L_\zeta)^{\frac{1}{2}} E_{\mathcal{D}}(\bar{u}).$$

- (Slow diffusion PME) Let  $\zeta(s) = |s|^{m-1}s$  with  $m \geq 1$ . There exists  $C_m > 0$  depending only on  $m$  such that

$$(43) \quad \left( \sum_{n=0}^{N-1} \delta t^{(n+\frac{1}{2})} \|\Pi_{\mathcal{D}} u^{(n+1)} - \Pi_{\mathcal{D}} I_{\mathcal{D}}^{n+1} \bar{u}\|_{L^{\frac{m+1}{m}}(\Omega)}^{m+1} \right)^{\frac{1}{m+1}} \leq C_m T^{\frac{1}{m+1}} E_{\mathcal{D}}(\bar{u})^{\frac{2}{m+1}}.$$

- (Fast diffusion PME) Let  $\zeta(s) = |s|^{m-1}s$  with  $m < 1$ . Then there exists  $C_m > 0$  depending only on  $m$  such that

$$(44) \quad \left( \sum_{n=0}^{N-1} \delta t^{(n+\frac{1}{2})} \|\zeta(\Pi_{\mathcal{D}} u^{(n+1)}) - \zeta(\Pi_{\mathcal{D}} I_{\mathcal{D}}^{n+1} \bar{u})\|_{L^{\frac{m+1}{m}}(\Omega)}^{\frac{m+1}{m}} \right)^{\frac{m}{m+1}} \leq C_m T^{\frac{m}{m+1}} E_{\mathcal{D}}(\bar{u})^{\frac{2m}{m+1}}.$$

*Corollary 3.7. Stefan equation.* Since  $\zeta$  is a Lipschitz-continuous and non-decreasing function, we have

$$|\zeta(a) - \zeta(b)|^2 \leq |\zeta(a) - \zeta(b)| L_\zeta |a - b| = L_\zeta (\zeta(a) - \zeta(b))(a - b), \quad \forall a, b \in \mathbb{R}.$$

Used in (39), this relation proves (42).

*Slow diffusion PME.* We first prove that, for some  $k_m > 0$ ,

$$(45) \quad |a - b|^{m+1} \leq k_m (|a|^{m-1}a - |b|^{m-1}b)(a - b), \quad \forall a, b \in \mathbb{R}.$$

The case for  $b = 0$  is trivial and reduces to  $k_m \geq 1$ . Consider  $b \neq 0$  and set  $s = a/b$ . To establish (45), we have to prove that  $|s - 1|^{m+1} \leq k_m (|s|^{m-1}s - 1)(s - 1)$ , which reduces to  $|s - 1|^m \leq c_m |s|^{m-1}s - 1|$ . The function  $s \mapsto \frac{|s-1|^m}{|s|^{m-1}s-1}$  is continuous on  $\mathbb{R}$  (use a Taylor expansion about  $s = 1$  to deal with the singularity) and has limit 1 at  $\pm\infty$ . It is therefore bounded, which proves the required estimate.

Using (45) in (39), the estimate (43) follows.

*Fast diffusion PME.* Let  $a', b' \in \mathbb{R}$  and apply (45) with  $\frac{1}{m} > 1$  instead of  $m$  and  $a = \zeta(a') = |a'|^{m-1}a'$ ,  $b = \zeta(b') = |b'|^{m-1}b'$ . Noting that  $|a|^{\frac{1}{m}-1}a = a'$  and  $|b|^{\frac{1}{m}-1}b = b'$ , we infer

$$|\zeta(a') - \zeta(b')|^{\frac{1}{m+1}} \leq k_{1/m}(a' - b')(\zeta(a') - \zeta(b')).$$

Used in (39) this establishes (44).  $\square$

**3.4. Proof of Theorem 3.4.** We follow the approach of [18] which consists in identifying an error equation on the discrete solution and the interpolate of the continuous solution, and estimating a consistency error.

To identify the error equation we integrate (4) over the time intervals  $(t^{(j)}, t^{(j+1)})$  and use (30) to obtain for all  $j \in \{0, \dots, N-1\}$

$$(46) \quad \int_{\Omega} \Pi_{\mathcal{D}} \left( e_{\mathcal{D},u}^{(j+1)} - e_{\mathcal{D},u}^j \right) \Pi_{\mathcal{D}} v + \delta t^{(j+\frac{1}{2})} \int_{\Omega} \nabla_{\mathcal{D}} e_{\mathcal{D},\zeta}^{(j+1)} \cdot \nabla_{\mathcal{D}} v = \mathfrak{E}_{\mathcal{D}}^j(v),$$

and all  $v \in X_{\mathcal{D},0}$ , where the consistency error is defined by

$$(47) \quad \mathfrak{E}_{\mathcal{D}}^j(v) := \delta t^{(j+\frac{1}{2})} \int_{\Omega} f^{(j+1)} \Pi_{\mathcal{D}} v - \int_{\Omega} \left[ \Pi_{\mathcal{D}} I_{\mathcal{D}}^{j+1} \bar{u} - \Pi_{\mathcal{D}} I_{\mathcal{D}}^j \bar{u} \right] \Pi_{\mathcal{D}} v - \delta t^{(j+\frac{1}{2})} \int_{\Omega} \nabla_{\mathcal{D}} \zeta(I_{\mathcal{D}}^{j+1} \bar{u}) \cdot \nabla_{\mathcal{D}} v.$$

This is a linear form  $\mathfrak{E}_{\mathcal{D}}^n : X_{\mathcal{D},0} \rightarrow \mathbb{R}$ . Its boundedness is established in

**Lemma 3.8.** For all  $v \in X_{\mathcal{D},0}$  and  $n \in \{0, \dots, N-1\}$ , there holds

$$(48) \quad |\mathfrak{E}_{\mathcal{D}}^n(v)| \leq \delta t^{(n+\frac{1}{2})} E_{\mathcal{D}}^n(\bar{u}) \|\nabla_{\mathcal{D}} v\|.$$

*Proof.* Recalling (34),  $\Pi_{\mathcal{D}} I_{\mathcal{D}}^{n+1} \bar{u} - \Pi_{\mathcal{D}} I_{\mathcal{D}}^n \bar{u}$  can be replaced by  $\bar{u}(t^{(n+1)}) - \bar{u}(t^{(n)})$  in (47) with a cost measured by  $S_{\mathcal{D}}^{\Pi,*}$ . Specifically, by the definition of the discrete  $H^{-1}$ -seminorm in (33) we have

$$(49) \quad \left| \int_{\Omega} \phi \Pi_{\mathcal{D}} v \right| \leq |\phi|_{\mathcal{D},*} \|\nabla_{\mathcal{D}} v\| \quad \forall \phi \in L^2(\Omega), \quad \forall v \in X_{\mathcal{D},0},$$

and thus, since  $I_{\mathcal{D}}^k \bar{u} = I_{\mathcal{D}}(\bar{u}(t^{(k)}))$  for  $k = n, n+1$ ,

$$\begin{aligned} & \left| \int_{\Omega} [\Pi_{\mathcal{D}} I_{\mathcal{D}}^{n+1} \bar{u} - \Pi_{\mathcal{D}} I_{\mathcal{D}}^n \bar{u}] \Pi_{\mathcal{D}} v - \int_{\Omega} [\bar{u}(t^{(n+1)}) - \bar{u}(t^{(n)})] \Pi_{\mathcal{D}} v \right| \\ &= \left| \int_{\Omega} [\Pi_{\mathcal{D}} I_{\mathcal{D}} (\bar{u}(t^{(n+1)}) - \bar{u}(t^{(n)})) - (\bar{u}(t^{(n+1)}) - \bar{u}(t^{(n)}))] \Pi_{\mathcal{D}} v \right| \\ &\leq S_{\mathcal{D}}^{\Pi,*}(\bar{u}(t^{(n+1)}) - \bar{u}(t^{(n)})) \|\nabla_{\mathcal{D}} v\|. \end{aligned}$$

Similarly, replacing  $\nabla_{\mathcal{D}} \zeta(I_{\mathcal{D}}^{n+1} \bar{u}) = \nabla_{\mathcal{D}} I_{\mathcal{D}} \zeta(\bar{u}(t^{(n+1)}))$  (see (29)) with  $\nabla \zeta(\bar{u}(t^{(n+1)}))$  incurs a cost measured by  $S_{\mathcal{D}}^{\nabla}(\zeta(\bar{u}(t^{(n+1)})))$ :

$$\left| \int_{\Omega} \nabla_{\mathcal{D}} \zeta(I_{\mathcal{D}}^{n+1} \bar{u}) \cdot \nabla_{\mathcal{D}} v - \int_{\Omega} \nabla \zeta(\bar{u}(t^{(n+1)})) \cdot \nabla_{\mathcal{D}} v \right| \leq S_{\mathcal{D}}^{\nabla}(\zeta(\bar{u}(t^{(n+1)}))) \|\nabla_{\mathcal{D}} v\|.$$

Hence,

$$\begin{aligned} \mathfrak{E}_{\mathcal{D}}^n(v) &= \delta t^{(n+\frac{1}{2})} \int_{\Omega} f^{(n+1)} \Pi_{\mathcal{D}} v - \int_{\Omega} [\bar{u}(t^{(n+1)}) - \bar{u}(t^{(n)})] \Pi_{\mathcal{D}} v - \delta t^{(n+\frac{1}{2})} \int_{\Omega} \nabla \zeta(\bar{u}(t^{(n+1)})) \cdot \nabla_{\mathcal{D}} v \\ &\quad + \mathcal{O}_1 \left[ S_{\mathcal{D}}^{\Pi,*}(\bar{u}(t^{(n+1)}) - \bar{u}(t^{(n)})) + \delta t^{(n+\frac{1}{2})} S_{\mathcal{D}}^{\nabla}(\zeta(\bar{u}(t^{(n+1)}))) \right] \|\nabla_{\mathcal{D}} v\| \end{aligned}$$

where, here and in the following,  $\mathcal{O}_1(X)$  denotes a generic function such that  $|\mathcal{O}_1(X)| \leq |X|$ . Now, by definition of  $W_{\mathcal{D}}(\nabla \zeta(\bar{u}(t^{(n+1)})))$ ,

$$\left| \int_{\Omega} \nabla \zeta(\bar{u}(t^{(n+1)})) \cdot \nabla_{\mathcal{D}} v + \Delta \zeta(\bar{u}(t^{(n+1)})) \Pi_{\mathcal{D}} v \right| \leq W_{\mathcal{D}}(\nabla \zeta(\bar{u}(t^{(n+1)}))) \|\nabla_{\mathcal{D}} v\|$$

and thus, writing  $\bar{u}(t^{(n+1)}) - \bar{u}(t^{(n)}) = \int_{t^{(n)}}^{t^{(n+1)}} \partial_t \bar{u}(s) ds$  (which is valid since the equation (1) and the regularity  $\nabla \zeta(\bar{u}) \in C([0, T]; H_{\text{div}}(\Omega))$  imply  $\partial_t \bar{u} \in C([0, T]; L^2(\Omega))$ ) and recalling the definition of  $f^{(n+1)}$ ,

$$\begin{aligned} (50) \quad \mathfrak{E}_{\mathcal{D}}^n(v) &= \int_{\Omega} \left[ \int_{t^{(n)}}^{t^{(n+1)}} (f - \partial_t \bar{u})(s) ds + \delta t^{(n+\frac{1}{2})} \Delta \zeta(\bar{u}(t^{(n+1)})) \right] \Pi_{\mathcal{D}} v \\ &\quad + \mathcal{O}_1 \left[ S_{\mathcal{D}}^{\Pi,*}(\bar{u}(t^{(n+1)}) - \bar{u}(t^{(n)})) + \delta t^{(n+\frac{1}{2})} S_{\mathcal{D}}^{\nabla}(\zeta(\bar{u}(t^{(n+1)}))) \right. \\ &\quad \left. + \delta t^{(n+\frac{1}{2})} W_{\mathcal{D}}(\nabla \zeta(\bar{u}(t^{(n+1)}))) \right] \|\nabla_{\mathcal{D}} v\|. \end{aligned}$$

Since  $f - \partial_t \bar{u} = -\Delta \zeta(\bar{u})$ , the property (49) yields

$$\begin{aligned} & \left| \int_{\Omega} \left[ \int_{t^{(n)}}^{t^{(n+1)}} (f - \partial_t \bar{u})(s) ds + \delta t^{(n+\frac{1}{2})} \Delta \zeta(\bar{u}(t^{(n+1)})) \right] \Pi_{\mathcal{D}} v \right| \\ &\leq \delta t^{(n+\frac{1}{2})} \left| \frac{1}{\delta t^{(n+\frac{1}{2})}} \int_{t^{(n)}}^{t^{(n+1)}} \Delta \zeta(\bar{u}(s)) ds - \Delta \zeta(\bar{u}(t^{(n+1)})) \right|_{\mathcal{D},*} \|\nabla_{\mathcal{D}} v\|. \end{aligned}$$

Plugging this estimate into (50) and recalling the definition (41) of  $E_{\mathcal{D}}^n(v)$ , this shows (48).  $\square$

With Lemma 3.8 at hand, the next proposition is the main step towards the error estimate in the fully discrete setting.

**Proposition 3.9.** *For all  $n \in \{1, \dots, N\}$ , there holds*

$$(51) \quad \max_{1 \leq n \leq N} \|\nabla_{\mathcal{D}} e_{\mathcal{D}, \zeta}^n\|^2 + 4 \sum_{n=0}^{N-1} \delta t^{(n+\frac{1}{2})} (\Pi_{\mathcal{D}} e_{\mathcal{D}, u}^{(n+1)}, \Pi_{\mathcal{D}} e_{\mathcal{D}, \zeta}^{(n+1)}) \leq 24T \exp(1) E_{\mathcal{D}}(\bar{u}).$$

*Proof.* Following the lines of the proof in the semi-discrete case, we sum (46) over  $j \in \{0, \dots, p-1\}$  with  $p \in \{1, \dots, N\}$ , leading to

$$(52) \quad \int_{\Omega} \Pi_{\mathcal{D}} e_{\mathcal{D}, u}^{(p)} \Pi_{\mathcal{D}} v + \sum_{j=0}^{p-1} \delta t^{(j+\frac{1}{2})} \int_{\Omega} \nabla_{\mathcal{D}} e_{\mathcal{D}, \zeta}^{(j+1)} \cdot \nabla_{\mathcal{D}} v = \sum_{j=0}^{p-1} \mathfrak{E}_{\mathcal{D}}^j(v).$$

Here, we used the fact that  $e_{\mathcal{D},u}^0 = 0$  thanks to the definition of  $u^0 = I_{\mathcal{D}}u_{\text{ini}}$ . Choosing  $v = \delta t^{(p-\frac{1}{2})}e_{\mathcal{D},\zeta}^{(p)}$  in the above equation and summing over  $p \in \{1, \dots, n\}$  for some  $n \in \{1, \dots, N\}$  provides

$$(53) \quad \mathfrak{J}_1^{(n)} + \mathfrak{J}_2^{(n)} = \mathfrak{R}^{(n)},$$

where

$$\begin{aligned} \mathfrak{J}_1^{(n)} &:= \sum_{p=1}^n \delta t^{(p-\frac{1}{2})} \int_{\Omega} \Pi_{\mathcal{D}} e_{\mathcal{D},u}^{(p)} \Pi_{\mathcal{D}} e_{\mathcal{D},\zeta}^{(p)}, \\ \mathfrak{J}_2^{(n)} &:= \sum_{p=1}^n \sum_{j=1}^p \delta t^{(p-\frac{1}{2})} \delta t^{(j-\frac{1}{2})} \int_{\Omega} \nabla_{\mathcal{D}} e_{\mathcal{D},\zeta}^{(j)} \cdot \nabla_{\mathcal{D}} e_{\mathcal{D},\zeta}^{(p)}, \\ \mathfrak{R}^{(n)} &:= \sum_{p=1}^n \sum_{j=0}^{p-1} \delta t^{(p-\frac{1}{2})} \mathfrak{E}_{\mathcal{D}}^j(e_{\mathcal{D},\zeta}^{(p)}). \end{aligned}$$

$\mathfrak{J}_1^{(n)}$  corresponds to the third term in (39). On the other hand, the identity (9) ensures that

$$(54) \quad \mathfrak{J}_2^{(n)} \geq \frac{1}{2} \int_{\Omega} \left| \nabla_{\mathcal{D}} \varepsilon_{\mathcal{D},\zeta}^{(n)} \right|^2.$$

The term  $\mathfrak{R}^{(n)}$  can be reorganized as

$$\mathfrak{R}^{(n)} = \sum_{j=0}^{n-1} \mathfrak{E}_{\mathcal{D}}^j(\varepsilon_{\mathcal{D},\zeta}^{(n)}) - \sum_{j=1}^{n-1} \mathfrak{E}_{\mathcal{D}}^j(\varepsilon_{\mathcal{D},\zeta}^{(j)}) := \mathfrak{R}_1^{(n)} + \mathfrak{R}_2^{(n)}.$$

Owing to Lemma 3.8, the first contribution  $\mathfrak{R}_1^{(n)}$  can be estimated as follows:

$$\left| \mathfrak{R}_1^{(n)} \right| \leq \sum_{j=0}^{n-1} \left| \mathfrak{E}_{\mathcal{D}}^j(\varepsilon_{\mathcal{D},\zeta}^{(n)}) \right| \leq \sum_{j=0}^{n-1} \delta t^{(j+\frac{1}{2})} E_{\mathcal{D}}^j(\bar{u}) \|\nabla_{\mathcal{D}} \varepsilon_{\mathcal{D},\zeta}^{(n)}\|.$$

Applying Cauchy–Schwarz inequality and then the Young inequality (12) provides

$$\left| \mathfrak{R}_1^{(n)} \right| \leq T \sum_{j=0}^{n-1} \delta t^{(j+\frac{1}{2})} E_{\mathcal{D}}^j(\bar{u})^2 + \frac{1}{4} \|\nabla_{\mathcal{D}} \varepsilon_{\mathcal{D},\zeta}^{(n)}\|^2 \leq T E_{\mathcal{D}}(\bar{u})^2 + \frac{1}{4} \|\nabla_{\mathcal{D}} \varepsilon_{\mathcal{D},\zeta}^{(n)}\|^2,$$

so that, in view of (54), there holds

$$(55) \quad \mathfrak{J}_2^{(n)} - \mathfrak{R}_1^{(n)} \geq \frac{1}{4} \|\nabla_{\mathcal{D}} \varepsilon_{\mathcal{D},\zeta}^{(n)}\|^2 - T E_{\mathcal{D}}(\bar{u})^2.$$

Using once again Lemma 3.8, one can bound the term  $\mathfrak{R}_2^{(n)}$  by

$$\mathfrak{R}_2^{(n)} \leq \sum_{j=1}^{n-1} \left| \mathfrak{E}_{\mathcal{D}}^j(\varepsilon_{\mathcal{D},\zeta}^{(j)}) \right| \leq \sum_{j=1}^{n-1} \delta t^{(j+\frac{1}{2})} E_{\mathcal{D}}^j(\bar{u}) \|\nabla_{\mathcal{D}} \varepsilon_{\mathcal{D},\zeta}^{(j)}\|.$$

Using the Cauchy–Schwarz and Young inequality as above then yields

$$(56) \quad \mathfrak{R}_2^{(n)} \leq 2T E_{\mathcal{D}}(\bar{u})^2 + \frac{1}{8T} \sum_{j=1}^{n-1} \delta t^{(j+\frac{1}{2})} \|\nabla_{\mathcal{D}} \varepsilon_{\mathcal{D},\zeta}^{(j)}\|^2.$$

Combining (55)–(56) in (53) provides

$$\|\nabla_{\mathcal{D}} \varepsilon_{\mathcal{D},\zeta}^{(n)}\|^2 + 4 \sum_{j=0}^{n-1} \delta t^{(j+\frac{1}{2})} (\Pi_{\mathcal{D}} e_{\mathcal{D},u}^{(j+1)}, \Pi_{\mathcal{D}} e_{\mathcal{D},\zeta}^{(j+1)}) \leq \frac{1}{2T} \sum_{j=1}^{n-1} \delta t^{(j+\frac{1}{2})} \|\nabla_{\mathcal{D}} \varepsilon_{\mathcal{D},\zeta}^{(j)}\|^2 + 12T E_{\mathcal{D}}(\bar{u})^2.$$

The generalized Gronwall Lemma [32, Lemma 5.1] then yields

$$\|\nabla_{\mathcal{D}} \varepsilon_{\mathcal{D},\zeta}^{(n)}\|^2 + 4 \sum_{j=0}^{n-1} \delta t^{(j+\frac{1}{2})} (\Pi_{\mathcal{D}} e_{\mathcal{D},u}^{(j+1)}, \Pi_{\mathcal{D}} e_{\mathcal{D},\zeta}^{(j+1)}) \leq 12T \exp(1) E_{\mathcal{D}}(\bar{u}).$$

Choosing  $n = N$ , and then  $n = \operatorname{argmax}_j \|\nabla_{\mathcal{D}} \varepsilon_{\mathcal{D},\zeta}^{(j)}\|^2$ , we obtain (51).  $\square$

With Proposition 3.9 we have estimated  $\Pi_{\mathcal{D}} e_{\mathcal{D},u}^{(n)} \Pi_{\mathcal{D}} e_{\mathcal{D},\zeta}^{(n)}$  and  $\nabla_{\mathcal{D}} \varepsilon_{\mathcal{D},\zeta}^{(n)}$ . The proof of Theorem 3.4 is concluded by establishing the discrete  $L^\infty(0, T; H^{-1}(\Omega))$  estimate on  $(\Pi_{\mathcal{D}} e_{\mathcal{D},u}^{(n)})_{1 \leq n \leq N}$ . This is obtained in the lemma below. Together with Proposition 3.9, this completes the proof of Theorem 3.4.

**Lemma 3.10.** *For all  $n \in \{1, \dots, N\}$ , there holds*

$$\left| \Pi_{\mathcal{D}} e_{\mathcal{D},u}^{(n)} \right|_{\mathcal{D},*}^2 \leq 2 \|\nabla_{\mathcal{D}} \varepsilon_{\mathcal{D},\zeta}^{(n)}\|^2 + 2T E_{\mathcal{D}}(\bar{u})^2.$$

*Proof.* Applying (52) for  $p = n$  and with  $v \in X_{\mathcal{D},0}$  such that  $\|\nabla_{\mathcal{D}} v\| \leq 1$ , using the Cauchy–Schwarz inequality and recalling (48), we have

$$\int_{\Omega} \Pi_{\mathcal{D}} e_{\mathcal{D},u}^{(n)} \Pi_{\mathcal{D}} v \leq \|\nabla_{\mathcal{D}} \varepsilon_{\mathcal{D},\zeta}^{(n)}\| + \sum_{j=0}^{n-1} \delta t^{(j+\frac{1}{2})} E_{\mathcal{D}}^j(\bar{u}).$$

Taking the supremum over such  $v$  gives

$$\|\Pi_{\mathcal{D}} e_{\mathcal{D},u}^{(n)}\|_{\mathcal{D},*} \leq \|\nabla_{\mathcal{D}} \varepsilon_{\mathcal{D},\zeta}^{(n)}\| + \sum_{j=0}^{n-1} \delta t^{(j+\frac{1}{2})} E_{\mathcal{D}}^j(\bar{u}) \leq \|\nabla_{\mathcal{D}} \varepsilon_{\mathcal{D},\zeta}^{(n)}\| + \sqrt{T} E_{\mathcal{D}}(\bar{u}),$$

and the proof follows straightforwardly.  $\square$

#### 4. NUMERICAL EXAMPLES

**4.1. Numerical results.** For the numerical tests, we consider the porous medium, equation in dimension 2, corresponding to (1) with  $\zeta(\bar{u}) = |\bar{u}|^{m-1} \bar{u}$  for  $m \in \{2, 3, 4\}$ . The computational domain is given by  $T = 1$  and  $\Omega = (0, 1)^2$ , and the exact solution is  $\bar{u}(t, x) = \mathcal{B}(t_0 + t, x - x_0)$ , where  $t_0 = 0.1$ ,  $x_0 = (0.5, 0.5)$  and

$$(57) \quad \mathcal{B}(t, x) = t^{-\frac{1}{m}} \left\{ \left[ C_{\mathcal{B}} - \frac{m-1}{4m^2} \left( \frac{|x|}{t^{\frac{1}{2m}}} \right)^2 \right]_+ \right\}^{\frac{1}{m-1}},$$

is the Barenblatt–Pattle solution. The initial solution is fixed by  $u_{\text{ini}} = \bar{u}(0, \cdot)$ . We choose  $C_{\mathcal{B}} = 0.005$ , so that  $\mathcal{B}$  remains equal to 0 on  $\partial\Omega$  during the entire simulation  $t \in [0, 1]$ . Note that by the offset  $t_0$  in  $\mathcal{B}$ , the singularity of this function at  $t = 0$  is avoided, and the initial condition satisfies Assumption (A2).

The simulations are run over three different mesh families: a family of (mostly) hexagonal meshes, a family of locally refined Cartesian meshes, and a family of triangular meshes. Examples of members of each family are provided in Figure 1. We consider uniform time steps. For the coarsest mesh in each family, the time step is  $\delta t^{(n+\frac{1}{2})} = 0.1$  for all  $n$ ; then, for each mesh refinement, the time step is divided by 4. Since we use implicit Euler time stepping, this means that the truncation error in time decay as  $\mathcal{O}(h^2)$ , where  $h$  is the mesh size; as our spatial methods (see below) are low order, in the best possible situation (linear equations, smooth exact solution), the maximal approximation rates are  $\mathcal{O}(h)$  on gradients and  $\mathcal{O}(h^2)$  on functions. The choice of time steps thus ensures that the spatial truncation error is the leading term in the estimate. The following schemes will be used for the tests.

- LEPNC (Locally Enriched Polytopal Non-Conforming scheme) [24]: applicable on generic polytopal meshes, one unknown per internal edge (after static condensation), based on broken polynomial functions with weak continuity properties across the edges. We have taken a zero weight  $\varpi$  on the edge unknowns, so  $\Pi_{\mathcal{D}}$  is only computed from the cell unknowns.
- HMM (Hybrid Mimetic Mixed scheme) [23, Chapter 13]: applicable on generic polytopal meshes, one unknown per internal edge (after static condensation), based on local reconstruction of piecewise constant functions and gradients.
- MLP1 (Mass-Lumped  $\mathbb{P}^1$  finite element) [23, Section 8.4]: only applicable on triangular meshes, one unknown per vertex, based on standard  $\mathbb{P}^1$  shape functions for the gradient and piecewise constant reconstruction around each vertex.
- VAG-a (Vertex Approximate Gradient, first presentation) [23, Section 8.5]: applicable on generic polytopal meshes, one unknown per internal vertex and one unknown per cell, based on standard  $\mathbb{P}^1$  on a triangular subdivision of the cells (using the center of the

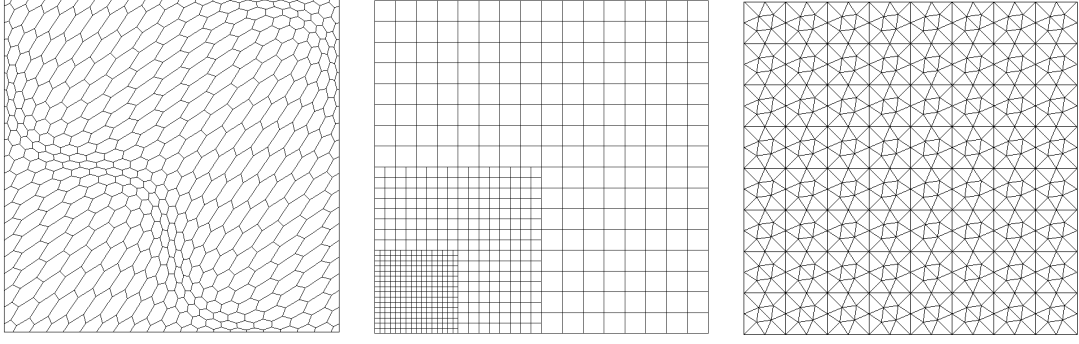


FIGURE 1. Examples of meshes used in numerical tests: hexagonal (left), locally refined Cartesian (centre), and triangular (right).

cell as additional vertex), with a mass-lumping that equally distributes the available area between cell and vertex unknowns. A local algebraic elimination (static condensation) of cell unknowns is also performed, leading to a globally coupled system on the vertex unknowns only.

- VAG-b (Vertex Approximate Gradient, second presentation) [13]: as above, but applied after writing the diffusion term as  $\text{div}(m|u|^{m-1}\nabla u)$ . Note that this scheme does not present itself as a gradient scheme.
- CFVEM (Conforming Virtual Element) [7]: applicable on generic polytopal meshes, one unknown per internal vertex, based on the elliptic projection of virtual shape functions, with algebraic mass-lumping.
- HDG (Hybridizable Discontinuous Galerkin, order 1) [17]: applicable on generic polytopal meshes (but the results are presented here only on triangular meshes), based on modal Legendre–Dubiner basis functions with one polynomial of degree 1 per cell and edge. The degrees of freedom are reduced to only edge polynomials after static condensation.

The LEPNC and HMM tests are based on the code available from in the **HArDCore2D** library [1] (based on the implementation principles of Hybrid High-Order methods [19]), while MLP1 tests were conducted using the code at <https://github.com/jdroniou/matlab-PME>.

**Remark 4.1** (CFVEM and  $\mathbb{P}^1$  finite elements on triangular meshes). *On triangular meshes and for the standard Laplace problem, CFVEM coincides with the conforming  $\mathbb{P}^1$  finite element method. The results presented below however show different behaviour of CFVEM and MLP1; the main reason can be found in the mass-lumping strategy adopted for each method: for MLP1, a geometrical mass-lumping was used, allocating to each vertex a mass corresponding to 1/3 of the sum of the areas of the triangles it belongs to; for CFVEM, an algebraic mass-lumping was used, reducing the standard mass matrix to a diagonal one by summing all elements on each row.*

The accuracy of the schemes are provided through the following quantities, all measured at the final time:

- Relative error in  $L^2$ -norm between the (reconstructed) gradients of the approximation of  $\zeta(\bar{u})$  and of the interpolate of  $\zeta(\bar{u})$ :

$$(58) \quad E_{H^1, \zeta} = \frac{\|\nabla_{\mathcal{D}}(\zeta(u^N) - I_{\mathcal{D}}\zeta(\bar{u})(T, \cdot))\|}{\|\nabla_{\mathcal{D}}I_{\mathcal{D}}\zeta(\bar{u})(T, \cdot)\|}.$$

- Relative error in  $L^{m+1}$ -norm between the (reconstructed) functions of the approximate solution and of the interpolate of the exact solution  $\bar{u}$ :

$$(59) \quad E_{L^{m+1}} = \frac{\|\Pi_{\mathcal{D}}(u^N - I_{\mathcal{D}}\bar{u}(T, \cdot))\|_{L^{m+1}(\Omega)}}{\|\Pi_{\mathcal{D}}I_{\mathcal{D}}\bar{u}(T, \cdot)\|_{L^{m+1}(\Omega)}}.$$



- Fraction of negative mass over total mass:

$$(60) \quad \text{NMass} = \frac{\int_{\Omega} (\Pi_{\mathcal{D}} u^N)_-}{\int_{\Omega} |\Pi_{\mathcal{D}} u^N|},$$

where  $s_- = \max(-s, 0)$  is the negative part of  $s \in \mathbb{R}$ .

**4.1.1. Rates of convergence vs. mesh size.** We first present the relative errors versus the mesh sizes, for all considered schemes and mesh families. The outputs are given in loglog graphs in Figures 2, 3 and 4. In these figures, the chosen reference slopes correspond to an estimate of the overall behaviour of the schemes, drawn from the tables as well as computed rates of convergence from one mesh to the other. Combining Estimates (39), (43) and Remark 3.5, and considering averaged-in-time norms (which are less stringent than the final time norms (58) and (59)), we would expect for a smooth enough exact solution a rate of convergence  $\mathcal{O}(h^{\frac{2}{m+1}})$  in  $L^{m+1}$ -norm on  $\bar{u}$  and  $\mathcal{O}(h)$  in  $L^2$ -norm on (the time integral of)  $\nabla \zeta(\bar{u})$ .

As can be seen in the numerical results, for the considered tests, the theoretical rates for  $E_{L^{m+1}}$  are sub-optimal for the lower values of  $m$  but tend to be close to the observed results for  $m = 4$ . The general trend, seen in both the theoretical and numerical results, is that the rates deteriorate as  $m$  increases. Two reasons can be found for that: as  $m$  increases, the  $L^{m+1}$ -norm becomes more constrained, while the regularity of the exact solution  $\bar{u}$  decays (for example, it is  $H^1$  in space for  $m = 2$ , but no longer for  $m > 2$ ).

Interestingly, but not surprisingly, the rates for  $E_{H^1, \zeta}$  seem to resist a little bit better as  $m$  increases, although they appear to be slightly below 1 for  $m > 2$  (which is not surprising since, as mentioned, (58) is a more constraining norm than  $\|\nabla \varepsilon_{\mathcal{D}, \zeta}^{(n)}\|$ ). Although  $E_{H^1, \zeta}$  measures an approximation of the gradient, which can be expected to be of lower order than that of a function, it measures this in a norm that is independent of  $m$  and relates to  $\zeta(\bar{u})$ , a function that has better regularity properties than  $\bar{u}$  (for example, it is  $H^1$  in space irrespective of the value of  $m$ ).

Comparing the various schemes, they all seem to adopt similar rates of convergence for  $E_{H^1, \zeta}$  on hexagonal and locally refined Cartesian meshes; the differences mostly lie in the multiplicative constants, with the largest factor between these multiplicative constants of order 10. More variation in the rates is observed on these mesh families for  $E_{L^{m+1}}$ , which is probably due to the variation of  $m$  and reduced regularity of  $\bar{u}$ , as discussed above. The rates on triangular meshes seem to depend much more on the chosen scheme. Focusing on  $E_{H^1, \zeta}$ , which is a more stable measure, we see that MLP1 outperforms the other schemes, at least for  $m = 2, 3$ ; of course, the drawback of MLP1 is that it can only be applied on triangular meshes. The other outlier is HDG, whose rates are much lower than the other schemes; the reason for that might be found in the total number of degrees of freedom, after static condensation, which is lower for HDG than some other schemes (see discussion in Section 4.1.2), and which therefore prevents this scheme from achieving optimal rates with respect to the mesh size.

It can also be noticed that some schemes produce a better  $E_{H^1, \zeta}$  error than others, but that the “ranking” between the schemes can be reversed if we look at the error  $E_{L^{m+1}}$ .

**4.1.2. Algebraic complexity.** We now briefly discuss the performance of the schemes relative to their algebraic complexity, measured primarily here in terms of their number **NDOFs** of degrees of freedom. Focusing only on  $m = 4$  (the most severe case), and hexagonal and triangular grids, we plot in Figure 5 the energy error  $E_{H^1, \zeta}$  of each scheme versus its degrees of freedom.

A first remark is that this measure is slightly more favourable to HDG than in the previous section. Its rate remains a bit lower than the other schemes, at least for the considered meshes, but perhaps less so than when comparing the error versus the mesh size. We also notice that, on triangular meshes, the (mostly) vertex-centered methods VAG-a, VAG-b, and MLP1 outperform the other schemes, which is expected: on triangular meshes, vertex-based methods have much fewer degrees of freedom than edge-based methods such as LENCP or HMM. Curiously, using the same argument, we would expect CFVEM, which is also a vertex-based method, to behave better than it does; this could be explained by the different kind of mass-lumping applied for these two methods.

On hexagonal meshes, except for LEPNC, all schemes presented here have a comparable error vs. complexity. The advantage of vertex-based methods is less perceptible on such meshes than on

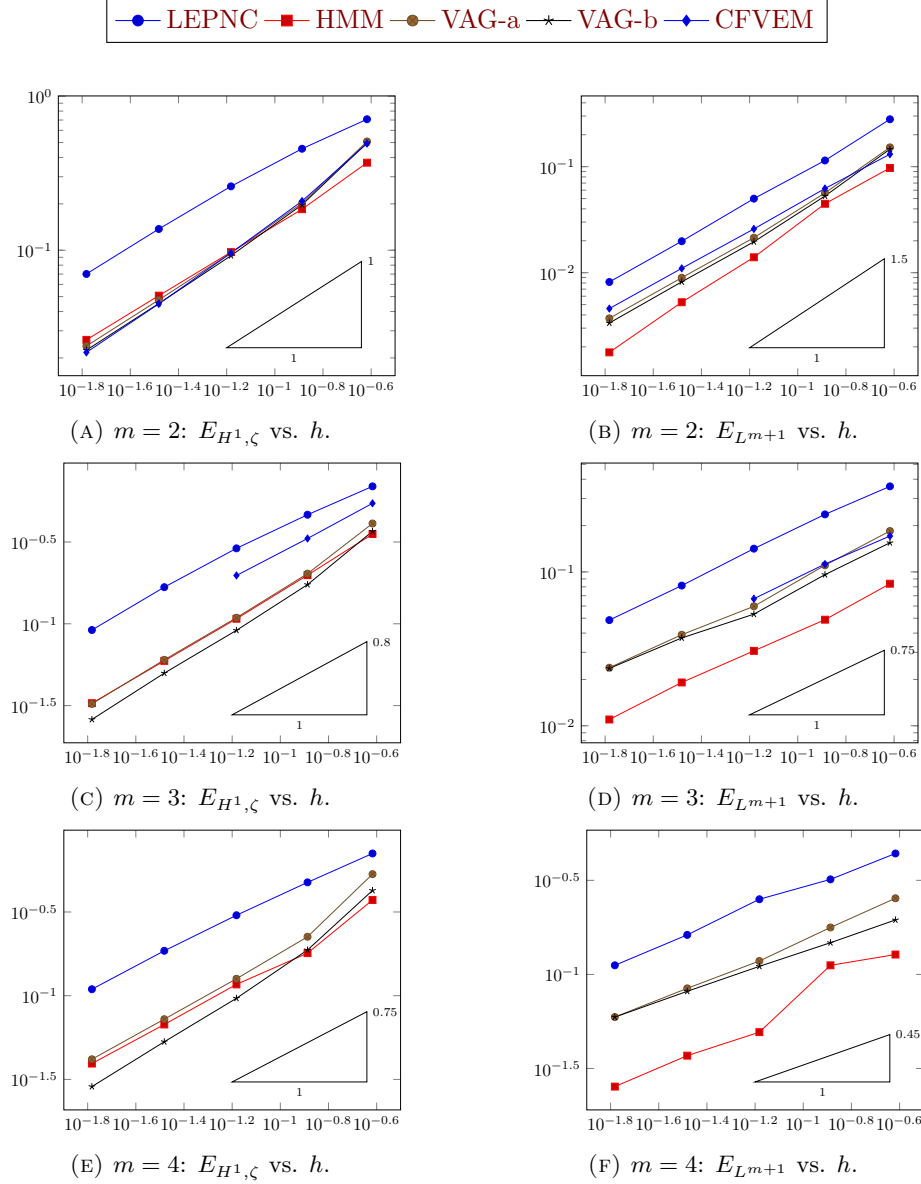


FIGURE 2. Hexagonal meshes: errors versus mesh size.

triangular meshes, which is not surprising: triangular meshes roughly have three times more edges than vertices, while hexagonal meshes have 1.5 times more edges than vertices on average.

**4.1.3. Positivity.** Finally, we look at the positivity properties of the schemes. As the standard linear heat equation, the porous medium equation satisfies a maximum principle: if the initial solution and the source terms are positive, then the solution remains positive for all times. Maintaining this property at the discrete level is particularly challenging, especially for schemes designed for generic polygonal/polyhedral meshes [20]. Except for MLP1 (which is restricted to triangular meshes), none of the schemes presented here satisfy this property in general.

In Figures 6, 7 and 8, we present the log-log graphs of the relative negative masses **NMass** versus the mesh sizes. In most situations, the schemes produce some negative mass, but it decays as the mesh is refined and is rather small relative to the total mass of the solution at the final time. On hexagonal and locally refined Cartesian meshes, the VAG and CFVEM schemes –which are (mostly) vertex-centered– have better positivity properties than the edge-centered schemes LEPNC and HMM, with VAG-b outperforming all the other schemes. Surprisingly, perhaps, VAG-a and VAG-b do not even produce negative values at the final time on locally refined meshes (and are thus absent from Figure 7). On triangular meshes, however, VAG-a produces more negative mass

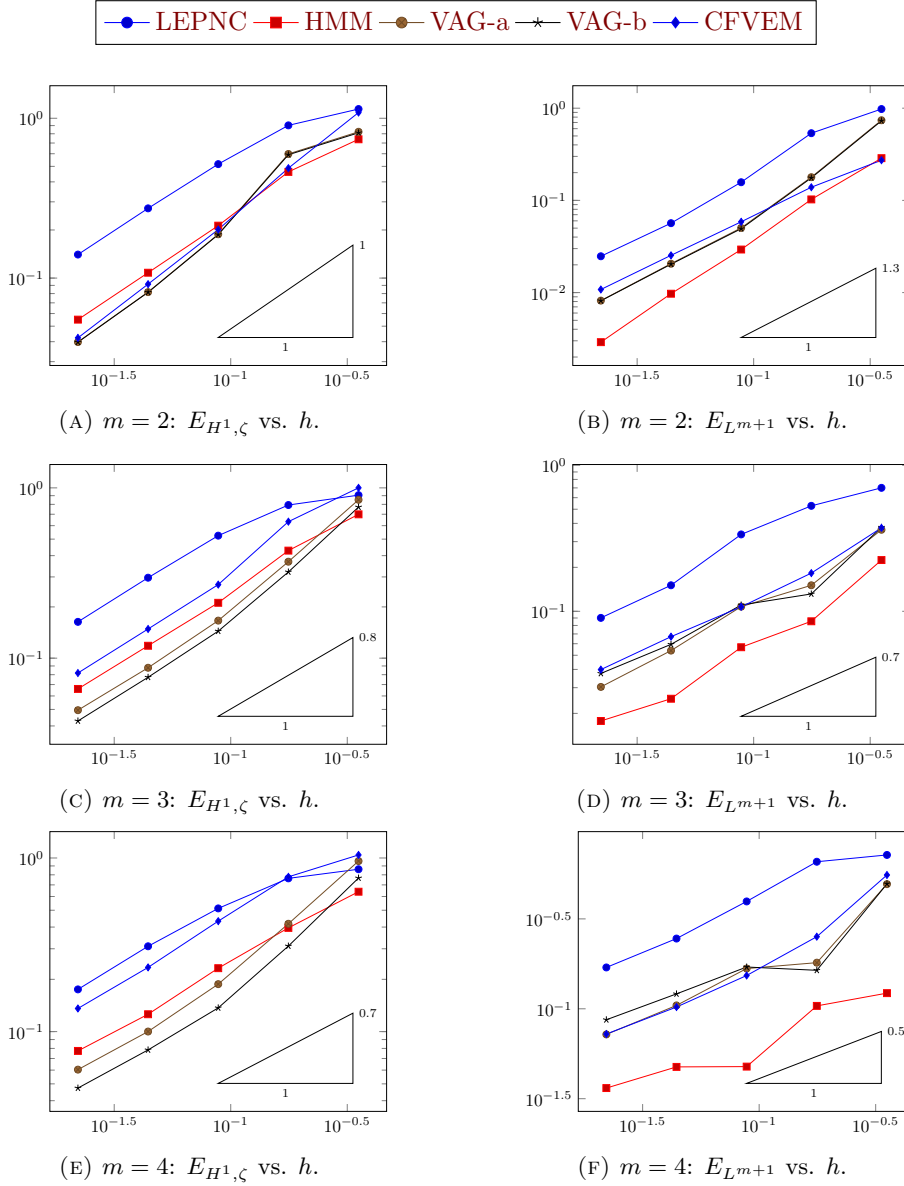


FIGURE 3. Locally refined Cartesian meshes: errors versus mesh size.

than the other schemes (and an amount comparable to HDG), with HMM and VAG-b being much closer to each other (with relative performance depending on  $m$ ), and LEPNC in between these clusters. MLP1 is known to preserve positive values and is therefore not represented in Figure 8. CFVEM was also found to preserve positive values on these triangular meshes.

These results demonstrate a strong interaction between scheme design and mesh geometries when it comes to preserving the maximum principle of the continuous model.

**4.2. A word on non-linear iterative schemes.** The time discrete problems (7) or the fully discrete counterparts (30) are nonlinear. To determine an approximation of the solution one needs to employ an iterative method. The common choice is the Newton method (see, e.g., [8]), which converges quadratically. However, this convergence is guaranteed only if the initial guess is close enough to the solution. A choice at hand being the solution computed at the previous time step, this means that the convergence is guaranteed if the differences between the solutions at two successive times are small enough. This induces restrictions on the time step.

For (1), the Newton method can fail to converge due to the singularities of  $\zeta$ , and in particular, for the fast diffusion case. To overcome this, one can regularise  $\zeta$  to avoid degeneracy, but even in this case, the convergence is only guaranteed under severe restrictions on the time step. To

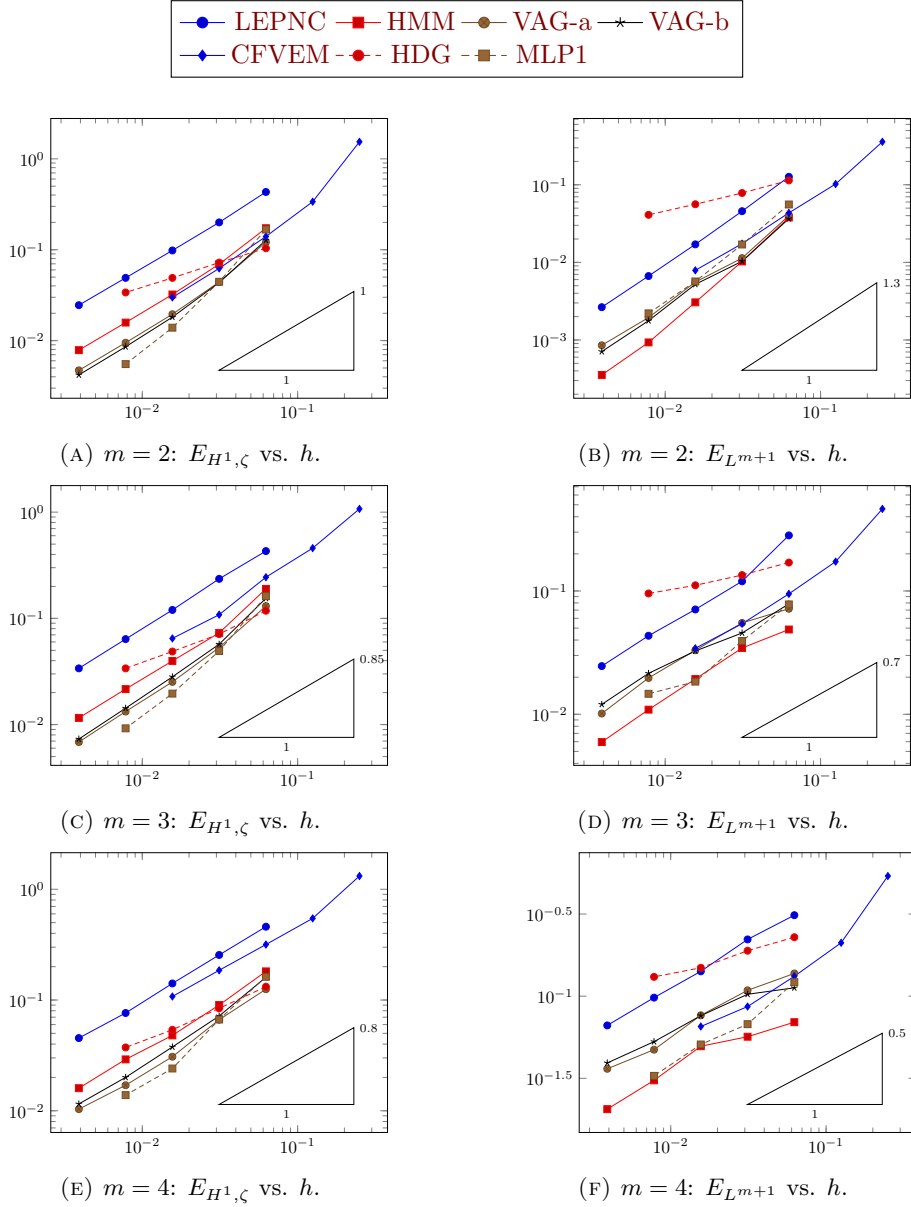


FIGURE 4. Triangular meshes: errors versus mesh size.

address these shortcomings, alternative iterative schemes have been designed. We mention here the relaxation scheme in [33], which shows to be more stable w.r.t. the choice of the initial condition, and the modified Picard scheme in [16], which is a simplified version of the Newton method. Both schemes are converging linearly. For these, as for the Newton scheme, the convergence is guaranteed rigorously under severe restrictions for the time step, as proved in [47].

A fixed point (contraction) scheme exploiting the monotonicity of  $\zeta$  has been proposed in [43] for the fast diffusion case and extended to more general situations in [46]. Though linear, the convergence is guaranteed under mild restrictions on the time step, regardless of the initial guess, and for any spatial discretisation. Moreover, as shown in [37], this scheme can be used to obtain a good initial guess for the Newton scheme, which leads to a stable and fast convergent iterative method. We also mention the scheme in [40], where the fixed point approach is combined with the Picard or Newton method by adding a stabilisation term. This leads to a scheme with the stability of the fixed point scheme and converging like the Picard scheme.

Finally, we refer to [9], where both  $\bar{u}$  and  $\zeta(\bar{u})$  are expressed in terms of a different unknown, based on a properly chosen parametrisation. This allows reformulating the Newton method in such

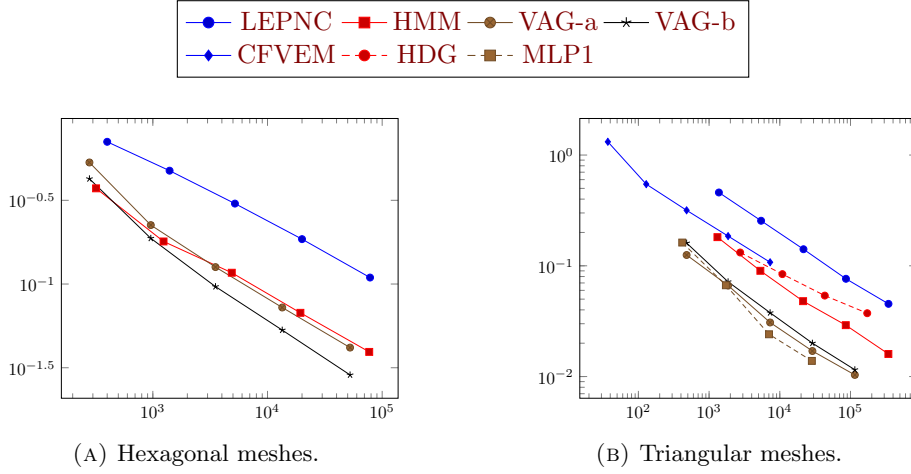
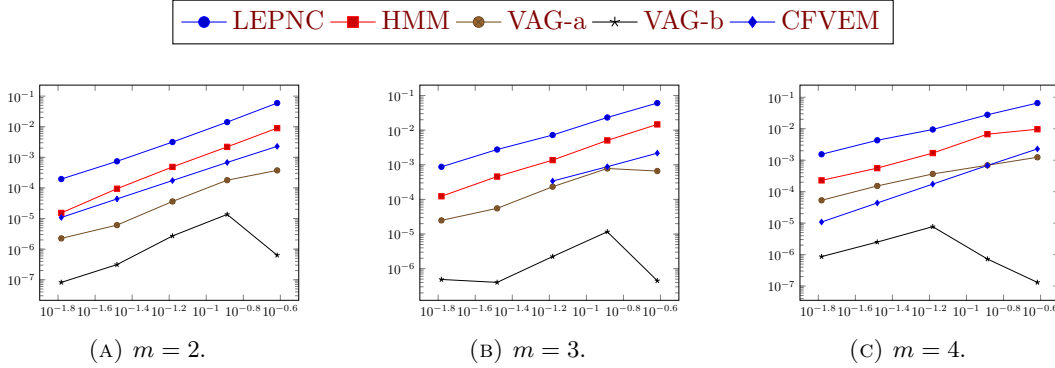
FIGURE 5. Energy error  $E_{H^1, \zeta}$  versus number of degrees of freedom NDOFs,  $m = 4$ .

FIGURE 6. Hexagonal meshes: fraction of negative mass NMass versus mesh size.

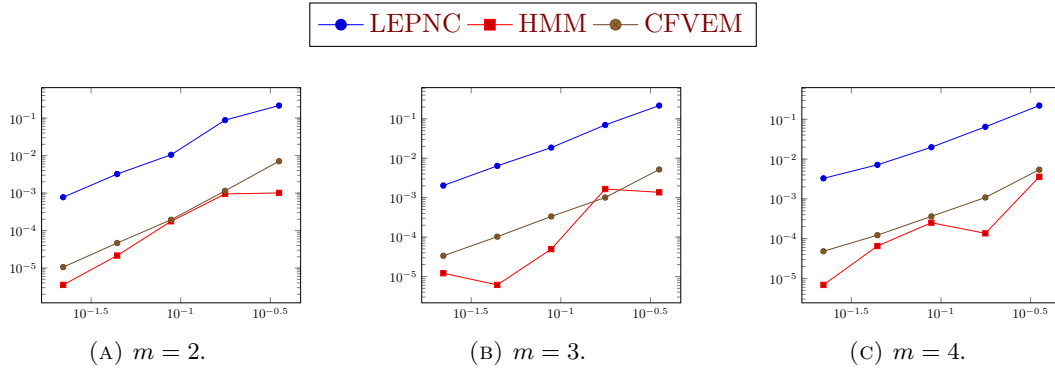
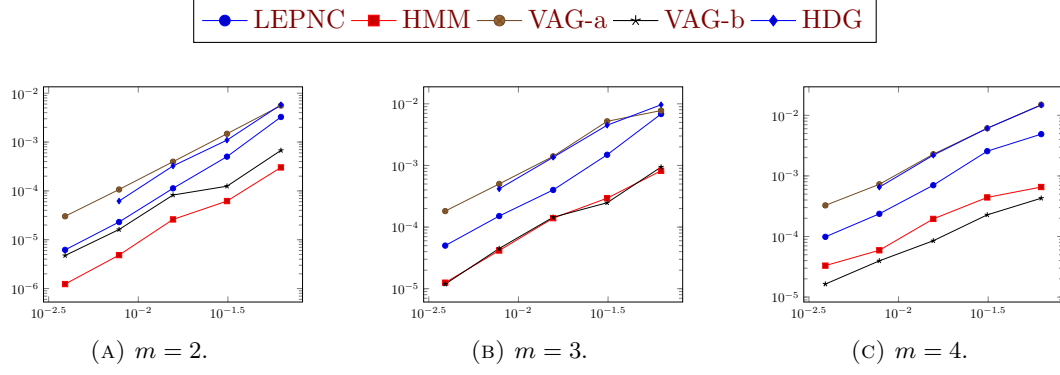


FIGURE 7. Locally refined Cartesian meshes: fraction of negative mass NMass versus mesh size.

a way that the quadratic convergence is unaffected, but the stability is significantly improved, so that larger time steps are allowed.

Here we discuss a simple iterative scheme that is inspired by the fixed point approach in [37, 43]. We restrict to the case where  $\zeta$  is Lipschitz continuous and let  $L_\zeta$  denote the Lipschitz constant, but the idea can be extended to more general situations as well, e.g., by applying the idea in [9]. For the ease of presentation, we present the scheme in the time-discrete case, the fully discrete one being analogous. With  $n \in \{0, \dots, N-1\}$  fixed we start by observing that (7) can be rewritten as

FIGURE 8. Triangular meshes: fraction of negative mass `NMass` versus mesh size.

the system

$$(61) \quad \begin{aligned} (\bar{u}^{(n+1)}, \varphi) + \delta t^{(n+\frac{1}{2})}(\nabla \bar{w}^{(n+1)}, \nabla \varphi) &= (\bar{u}^{(n)}, \varphi) + \delta t^{(n+\frac{1}{2})}(f^{(n+1)}, \varphi), \\ (\bar{w}^{(n+1)}, \psi) &= (\zeta(\bar{u}^{(n+1)}), \psi), \end{aligned}$$

for all  $\varphi \in H_0^1(\Omega)$  and  $\psi \in L^2(\Omega)$ . For a given  $L \geq \frac{L_\zeta}{2}$ , the iterative scheme consists in finding the pairs  $(\bar{u}^i, \bar{w}^i) \in H_0^1(\Omega) \times L^2(\Omega)$  ( $i \in \mathbb{N}^*$ ) solving the linear systems

$$(62) \quad \begin{aligned} (\bar{u}^i, \varphi) + \delta t^{(n+\frac{1}{2})}(\nabla \bar{w}^i, \nabla \varphi) &= (\bar{u}^{(n)}, \varphi) + \delta t^{(n+\frac{1}{2})}(f^{(n+1)}, \varphi), \\ (\bar{w}^i, \psi) &= L(\bar{u}^i - \bar{u}^{i-1}, \psi) + (\zeta(\bar{u}^{i-1}), \psi), \end{aligned}$$

for all  $\varphi \in H_0^1(\Omega)$  and  $\psi \in L^2(\Omega)$ . A natural choice for the initial guess is  $\bar{u}^0 = \bar{u}^{(n)}$  (the solution at the previous time step), but, as will be seen below, the convergence is guaranteed for any starting point. To prove this, we define the iteration errors

$$(63) \quad e_u^i = \bar{u}^{(n+1)} - \bar{u}^i, \quad \text{and} \quad e_w^i = \zeta(\bar{u}^{(n+1)}) - \bar{w}^i.$$

From (61) and (62), the errors satisfy

$$(64) \quad \begin{aligned} (e_u^i, \varphi) + \delta t^{(n+\frac{1}{2})}(\nabla e_w^i, \nabla \varphi) &= 0, \\ (e_w^i, \psi) &= L(e_u^i - e_u^{i-1}, \psi) + (\zeta(\bar{u}^{(n+1)}) - \zeta(\bar{u}^{i-1}), \psi), \end{aligned}$$

for all  $\varphi \in H_0^1(\Omega)$  and  $\psi \in L^2(\Omega)$ . With this, the convergence result is

**Lemma 4.2.** *The iterative scheme in (62) is convergent regardless of the initial guess. More precisely, one has  $\bar{w}^i \rightarrow \zeta(\bar{u}^{(n+1)})$  in  $H^1(\Omega)$  and  $\bar{u}^i \rightarrow \bar{u}^{(n+1)}$  in  $L^2(\Omega)$  as  $i \rightarrow \infty$ .*

*Proof.* Taking  $\varphi = e_w^i$  and  $\psi = e_u^i$  into (64) and subtracting the resulting gives

$$L\|e_u^i\|^2 + \delta t^{(n+\frac{1}{2})}\|\nabla e_w^i\|^2 = (Le_u^{i-1} - (\zeta(\bar{u}^{(n+1)}) - \zeta(\bar{u}^{i-1})), e_u^i).$$

Since  $\zeta$  is Lipschitz, by the choice of  $L$  one has  $|Le_u^{i-1} - (\zeta(\bar{u}^{(n+1)}) - \zeta(\bar{u}^{i-1}))| \leq L|e_u^{i-1}|$ . This, together with the Cauchy-Schwarz inequality leads to

$$L\|e_u^i\|^2 + \delta t^{(n+\frac{1}{2})}\|\nabla e_w^i\|^2 \leq L\|e_u^{i-1}\|\|e_u^i\|.$$

Applying now (12) and multiplying the resulting by 2 yields

$$(65) \quad L\|e_u^i\|^2 + 2\delta t^{(n+\frac{1}{2})}\|\nabla e_w^i\|^2 \leq L\|e_u^{i-1}\|^2.$$

Adding (65) for  $i = 1, \dots, k$  ( $k$  being arbitrary) leads to

$$(66) \quad L\|e_u^k\|^2 + 2\delta t^{(n+\frac{1}{2})} \sum_{i=1}^k \|\nabla e_w^i\|^2 \leq L\|e_u^0\|^2.$$

This shows that the second term above is a convergent series, implying the first convergence result. Using this convergence in the first equation of (64) completes the proof.  $\square$

## ACKNOWLEDGMENTS

CC acknowledges support from the Labex CEMPI (ANR-11-LABX-0007-01). JD and GM were partially supported by the Australian Government through the Australian Research Council's Discovery Projects funding scheme (project DP170100605); GM was also partially supported by the ERC Project CHANGE, which has received funding from the European Research Council (ERC) under the European Unions Horizon 2020 research and innovation programme (grant agreement No 694515). MBO and ISP are supported by the the Research Foundation-Flanders (FWO), Belgium through the Odysseus programme project G0G1316N.

## REFERENCES

- [1] HARDCore2D – Hybrid Arbitrary Degree::Core 2D. <https://github.com/jdroniou/HARDCore2D-release>, Version 2.0.2.
- [2] H. W. Alt and S. Luckhaus. Quasilinear elliptic-parabolic differential equations. *Math. Z.*, 183:311–341, 1983.
- [3] B. Andreianov, C. Cancès, and A. Moussa. A nonlinear time compactness result and applications to discretization of degenerate parabolic-elliptic PDEs. *J. Funct. Anal.*, 273(12):3633–3670, 2017.
- [4] O. Angelini, K. Brenner, and D. Hilhorst. A finite volume method on general meshes for a degenerate parabolic convection-reaction-diffusion equation. *Numer. Math.*, 123(2):219–257, 2013.
- [5] T. Arbogast, M. F. Wheeler, and N.-Y. Zhang. A nonlinear mixed finite element method for a degenerate parabolic equation arising in flow in porous media. *SIAM J. Numer. Anal.*, 33(4):1669–1687, 1996.
- [6] J. W. Barrett and P. Knabner. Finite element approximation of the transport of reactive solutes in porous media. II. Error estimates for equilibrium adsorption processes. *SIAM J. Numer. Anal.*, 34(2):455–479, 1997.
- [7] L. Beirão da Veiga, F. Brezzi, A. Cangiani, G. Manzini, L. D. Marini, and A. Russo. Basic principles of virtual element methods. *Math. Models Methods Appl. Sci. (M3AS)*, 199(23):199–214, 2013.
- [8] L. Bergamaschi and M. Putti. Mixed finite elements and Newton-type linearizations for the solution of Richards' equation. *Internat. J. Numer. Methods Engrg.*, 45(8):1025–1046, 1999.
- [9] K. Brenner and C. Cancès. Improving Newton's method performance by parametrization: the case of the Richards equation. *SIAM J. Numer. Anal.*, 55(4):1760–1785, 2017.
- [10] C. Cancès. Energy stable numerical methods for porous media flow type problems. *Oil Gas Sci. Technol. – Rev. IFP Energies nouvelles*, 73, 2018.
- [11] C. Cancès and T. Gallouët. On the time continuity of entropy solutions. *J. Evol. Equ.*, 11(1):43–55, 2011.
- [12] C. Cancès and C. Guichard. Convergence of a nonlinear entropy diminishing Control Volume Finite Element scheme for solving anisotropic degenerate parabolic equations. *Math. Comp.*, 85(298):549–580, 2016.
- [13] C. Cancès and C. Guichard. Numerical analysis of a robust free energy diminishing finite volume scheme for parabolic equations with gradient structure. *Found. Comput. Math.*, 17(6):1525–1584, 2017.
- [14] C. Cancès, I. S. Pop, and M. Vohralík. An a posteriori error estimate for vertex-centered finite volume discretizations of immiscible incompressible two-phase flow. *Math. Comp.*, 83(285):153–188, 2014.
- [15] J. Carrillo. Entropy solutions for nonlinear degenerate problems. *Arch. Ration. Mech. Anal.*, 147(4):269–361, 1999.
- [16] M. Celia, E. Bouloutas, and R. Zarba. A General Mass-Conservative Numerical-Soution for the Unsaturated Flow Equation. *Water Resour. Res.*, 26(7):1483–1496, 1990.
- [17] B. Cockburn, J. Gopalakrishnan, and R. Lazarov. Unified hybridization of discontinuous Galerkin, mixed, and continuous Galerkin methods for second order elliptic problems. *SIAM J. Numer. Anal.*, 47(2):1319–1365, 2009.
- [18] D. A. Di Pietro and J. Droniou. A third Strang lemma and an Aubin-Nitsche trick for schemes in fully discrete formulation. *Calcolo*, 55(3):Art. 40, 39p, 2018.
- [19] D. A. Di Pietro and J. Droniou. *The Hybrid High-Order Method for Polytopal Meshes: Design, Analysis, and Applications*, volume 19 of *Modeling, Simulation and Applications*. Springer International Publishing, 2020.
- [20] J. Droniou. Finite volume schemes for diffusion equations: introduction to and review of modern methods. *Math. Models Methods Appl. Sci. (M3AS)*, 24(8):1575–1619, 2014. Special issue on Recent Techniques for PDE Discretizations on Polyhedral Meshes.
- [21] J. Droniou and R. Eymard. Uniform-in-time convergence of numerical methods for non-linear degenerate parabolic equations. *Numer. Math.*, 132(4):721–766, 2016.
- [22] J. Droniou and R. Eymard. High-order mass-lumped schemes for nonlinear degenerate elliptic equations. *SIAM J. Numer. Anal.*, 58(1):153–188, 2020.
- [23] J. Droniou, R. Eymard, T. Gallouët, C. Guichard, and R. Herbin. *The gradient discretisation method*, volume 82 of *Mathematics & Applications*. Springer, 2018.
- [24] J. Droniou, R. Eymard, T. Gallouët, and R. Herbin. Non-conforming finite elements on polytopal meshes.
- [25] J. Droniou, R. Eymard, T. Gallouët, and R. Herbin. Gradient schemes: a generic framework for the discretisation of linear, nonlinear and nonlocal elliptic and parabolic equations. *Math. Models Methods Appl. Sci. (M3AS)*, 23(13):2395–2432, 2013.
- [26] Y. Epshteyn and B. Rivière. Analysis of  $hp$  discontinuous Galerkin methods for incompressible two-phase flow. *J. Comput. Appl. Math.*, 225(2):487–509, 2009.
- [27] A. Ern and I. Mozolevski. Discontinuous Galerkin method for two-component liquid-gas porous media flows. *Comput. Geosci.*, 16(3):677–690, 2012.
- [28] R. Eymard, P. Féron, T. Gallouët, C. Guichard, and R. Herbin. Gradient schemes for the Stefan problem. *Int. J. Finite Vol.*, 13:1–37, 2013.



- [29] R. Eymard, T. Gallouët, D. Hilhorst, and Y. Naït Slimane. Finite volumes and nonlinear diffusion equations. *RAIRO Modél. Math. Anal. Numér.*, 32(6):747–761, 1998.
- [30] R. Eymard, R. Herbin, and A. Michel. Mathematical study of a petroleum-engineering scheme. *M2AN Math. Model. Numer. Anal.*, 37(6):937–972, 2003.
- [31] R. Eymard, D. Hilhorst, and M. Vohralík. A combined finite volume–nonconforming/mixed-hybrid finite element scheme for degenerate parabolic problems. *Numer. Math.*, 105(1):73–131, 2006.
- [32] J. G. Heywood and R. Rannacher. Finite-element approximation of the nonstationary Navier-Stokes problem. IV. Error analysis for second-order time discretization. *SIAM J. Numer. Anal.*, 27(2):353–384, 1990.
- [33] W. Jäger and J. Kačur. Solution of doubly nonlinear and degenerate parabolic problems by relaxation schemes. *RAIRO Modél. Math. Anal. Numér.*, 29(5):605–627, 1995.
- [34] R. A. Klausen, F. A. Radu, and G. T. Eigestad. Convergence of MPFA on triangulations and for Richards’ equation. *Internat. J. Numer. Methods Fluids*, 58(12):1327–1351, 2008.
- [35] O. A. Ladyženskaja, V. A. Solonnikov, and N. N. Ural’ceva. *Linear and quasilinear equations of parabolic type*. Translated from the Russian by S. Smith. Translations of Mathematical Monographs, Vol. 23. American Mathematical Society, Providence, R.I., 1967.
- [36] N. Liao. A unified approach to the Hölder regularity of solutions to degenerate and singular parabolic equations. *J. Differential Equations*, 268(10):5704–5750, 2020.
- [37] F. List and F. A. Radu. A study on iterative methods for solving Richards’ equation. *Comput. Geosci.*, 20(2):341–353, 2016.
- [38] E. Magenes, R. H. Nochetto, and C. Verdi. Energy error estimates for a linear scheme to approximate nonlinear parabolic problems. *ESAIM: Mathematical Modelling and Numerical Analysis*, 21(4):655–678, 1987.
- [39] A. M. Meirmanov. *The Stefan problem*, volume 3 of *de Gruyter Expositions in Mathematics*. Walter de Gruyter & Co., Berlin, 1992. Translated from the Russian by Marek Niezgódka and Anna Crowley, With an appendix by the author and I. G. Götz.
- [40] K. Mitra and I. S. Pop. A modified L-scheme to solve nonlinear diffusion problems. *Comput. Math. Appl.*, 77(6):1722–1738, 2019.
- [41] R. H. Nochetto and C. Verdi. Approximation of degenerate parabolic problems using numerical integration. *SIAM J. Numer. Anal.*, 25(4):784–814, 1988.
- [42] F. Otto.  $L^1$ -contraction and uniqueness for quasilinear elliptic-parabolic equations. *J. Differ. Equations*, 131:20–38, 1996.
- [43] I. S. Pop, F. Radu, and P. Knabner. Mixed finite elements for the Richards’ equation: linearization procedure. *J. Comput. Appl. Math.*, 168(1-2):365–373, 2004.
- [44] I. S. Pop and B. Schweizer. Regularization schemes for degenerate Richards equations and outflow conditions. *Math. Models Methods Appl. Sci.*, 21(8):1685–1712, 2011.
- [45] I. S. Pop, M. Sepúlveda, F. A. Radu, and O. P. Vera Villagrán. Error estimates for the finite volume discretization for the porous medium equation. *J. Comput. Appl. Math.*, 234(7):2135–2142, 2010.
- [46] F. A. Radu, K. Kumar, J. M. Nordbotten, and I. S. Pop. A robust, mass conservative scheme for two-phase flow in porous media including Hölder continuous nonlinearities. *IMA J. Numer. Anal.*, 38(2):884–920, 2018.
- [47] F. A. Radu, I. S. Pop, and P. Knabner. Newton-type methods for the mixed finite element discretization of some degenerate parabolic equations. In *Numerical mathematics and advanced applications*, pages 1192–1200. Springer, Berlin, 2006.
- [48] F. A. Radu, I. S. Pop, and P. Knabner. Error estimates for a mixed finite element discretization of some degenerate parabolic equations. *Numer. Math.*, 109(2):285–311, 2008.
- [49] J. L. Vázquez. *Smoothing and decay estimates for nonlinear diffusion equations*, volume 33 of *Oxford Lecture Series in Mathematics and its Applications*. Oxford University Press, Oxford, 2006. Equations of porous medium type.
- [50] J. L. Vázquez. *The porous medium equation*. Oxford Mathematical Monographs. The Clarendon Press Oxford University Press, Oxford, 2007. Mathematical theory.
- [51] M. Vohralík and M. F. Wheeler. A posteriori error estimates, stopping criteria, and adaptivity for two-phase flows. *Comput. Geosci.*, 17(5):789–812, 2013.
- [52] C. S. Woodward and C. N. Dawson. Analysis of expanded mixed finite element methods for a nonlinear parabolic equation modeling flow into variably saturated porous media. *SIAM J. Numer. Anal.*, 37(3):701–724, 2000.
- [53] I. Yotov. A mixed finite element discretization on non-matching multiblock grids for a degenerate parabolic equation arising in porous media flow. *East-West J. Numer. Math.*, 5(3):211–230, 1997.
- [54] W. P. Ziemer. Interior and boundary continuity of weak solutions of degenerate parabolic equations. *Trans. Amer. Math. Soc.*, 271:733–748, 1982.

CLÉMENT CANCES ([clement.cances@inria.fr](mailto:clement.cances@inria.fr)): INRIA, UNIV. LILLE, CNRS, UMR 8524 - LABORATOIRE PAUL PAINLEVÉ, F-59000 LILLE.

JÉRÔME DRONIOU ([jerome.droniou@monash.edu](mailto:jerome.droniou@monash.edu)): SCHOOL OF MATHEMATICS, MONASH UNIVERSITY, MELBOURNE, AUSTRALIA.

CINDY GUICHARD ([cindy.guichard@sorbonne-universite.fr](mailto:cindy.guichard@sorbonne-universite.fr)): SORBONNE UNIVERSITÉ, CNRS, UNIVERSITÉ DE PARIS, INRIA ANGE, LABORATOIRE JACQUES-LOUIS LIONS, F-75005 PARIS, FRANCE.

GIANMARCO MANZINI ([marco.manzini@imati.cnr.it](mailto:marco.manzini@imati.cnr.it)): ISTITUTO DI MATEMATICA APPLICATA E TECNOLOGIE INFORMATICHE - CNR, VIA FERRATA 1, PAVIA, ITALY.

MANUELA BASTIDAS OLIVARES ([manuela.bastidas@uhasselt.be](mailto:manuela.bastidas@uhasselt.be)): HASSELT UNIVERSITY, CAMPUS DIEPENBEEK, AGORALAAN GEBOUW D, 3590 DIEPENBEEK, BELGIUM.

IULIU SORIN POP ([sorin.pop@uhasselt.be](mailto:sorin.pop@uhasselt.be)): HASSELT UNIVERSITY, CAMPUS DIEPENBEEK, AGORALAAN GEBOUW D, 3590 DIEPENBEEK, BELGIUM.